

Test Statistics Null Distributions in Multiple Testing: Simulation Studies and Applications to Genomics

Katherine S. Pollard*

Merrill D. Birkner[†]

Mark J. van der Laan[‡]

Sandrine Dudoit**

*Dept. of Biomolecular Engineering, University of California, Santa Cruz, kpollard@gladstone.ucsf.edu

[†]Division of Biostatistics, School of Public Health, University of California, Berkeley, mbirkner@berkeley.edu

[‡]Division of Biostatistics, School of Public Health, University of California, Berkeley, laan@berkeley.edu

**Division of Biostatistics, School of Public Health, University of California, Berkeley, sandrine@stat.berkeley.edu

This working paper is hosted by The Berkeley Electronic Press (bepress) and may not be commercially reproduced without the permission of the copyright holder.

<http://biostats.bepress.com/ucbbiostat/paper184>

Copyright ©2005 by the authors.

Test Statistics Null Distributions in Multiple Testing: Simulation Studies and Applications to Genomics

Katherine S. Pollard, Merrill D. Birkner, Mark J. van der Laan, and Sandrine Dudoit

Abstract

Multiple hypothesis testing problems arise frequently in biomedical and genomic research, for instance, when identifying differentially expressed or co-expressed genes in microarray experiments. We have developed generally applicable resampling-based single-step and stepwise multiple testing procedures (MTP) for control of a broad class of Type I error rates, defined as tail probabilities and expected values for arbitrary functions of the numbers of false positives and rejected hypotheses (Dudoit and van der Laan, 2005; Dudoit et al., 2004a,b; Pollard and van der Laan, 2004; van der Laan et al., 2005, 2004a,b). As argued in the early article of Pollard and van der Laan (2004), a key feature of the methodology is the general characterization and explicit construction of a test statistics null distribution (rather than data generating null distribution), which provides Type I error control in testing problems involving general data generating distributions (with arbitrary dependence structures among variables), null hypotheses, and test statistics. In particular, the proposed null distribution provides Type I error control without requirements such as subset pivotality (Westfall and Young, 1993) and, therefore, allows one to test hypotheses about a much broader class of parameters than are covered by currently available methods (e.g., correlation coefficients, regression parameters in linear and non-linear models with dependent covariates and error terms).

This paper presents simulation studies comparing test statistics null distributions in two testing scenarios of great relevance to biomedical and genomic data analysis: tests for regression parameters in linear models where covariates and error

terms are allowed to be dependent and tests for correlation coefficients. The simulation studies demonstrate that the choice of null distribution can have a substantial impact on the Type I error and power properties of a given multiple testing procedure. Procedures based on a general non-parametric bootstrap estimator of the proposed test statistics null distribution typically control the Type I error rate “on target” at the nominal level. In contrast, comparable procedures, based on parameter-specific bootstrap null distributions, can be severely anti-conservative (bootstrapping residuals for the test of regression parameters) or conservative (independent bootstrap for the test of correlation coefficients). Applications to a novel genomic dataset, from a study of microRNA expression in cancer, illustrate the flexibility and power of our proposed methodology (Lu et al., 2005).

1 Introduction

1.1 Motivation

The genomic age has brought growing interest in multiple testing. As new technologies, such as microarrays, mass spectrometry, and high-throughput sequencing, facilitate the collection of high-dimensional biological datasets, researchers are becoming increasingly reliant on statistical methods for assessing the significance of biological findings over families of thousands or even millions of hypothesis tests. Identifying differentially expressed or co-expressed genes from genome-wide mRNA abundance data are classic examples. Other applications include: tests of association between gene expression measures and Gene Ontology (GO) annotation (Dudoit and van der Laan (2005); www.geneontology.org); the identification of transcription factor binding sites in ChIP-Chip experiments, where chromatin immunoprecipitation (ChIP) of transcription factor-bound DNA is followed by microarray hybridization (Chip) of the IP-enriched DNA (Keleş et al., 2004); tests of association between phenotypes and amino acid mutations, e.g., viral replication capacity and HIV-1 sequence variation (Birkner et al., 2005b,c; van der Laan et al., 2005); the genetic mapping of complex traits using single nucleotide polymorphisms (SNP) (Birkner et al., 2005a). These testing problems are particularly challenging, as they involve inference for large multivariate distributions, with complex and unknown dependence structures among variables. Therefore, existing methods, based solely on the marginal distributions of the test statistics and/or simplifying assumptions about their joint distribution, are generally not appropriate.

Motivated by the aforementioned biomedical and genomic applications and the limitations of existing multiple testing methods, we have developed and implemented (in R and SAS) resampling-based single-step and stepwise *multiple testing procedures* (MTP) for controlling a broad class of Type I error rates, in testing problems involving general data generating distributions (with arbitrary dependence structures among variables), null hypotheses (defined in terms of submodels for the data generating distribution), and test statistics (e.g., t -statistics, F -statistics) (Dudoit and van der Laan, 2005; Dudoit et al., 2004a,b; Pollard and van der Laan, 2004; van der Laan et al., 2005, 2004a,b). In particular, procedures that take into account the *joint* distribution of the test statistics are provided to control Type I error rates defined as tail probabilities and expected values for arbitrary functions $g(V_n, R_n)$ of

the numbers of false positives V_n and rejected hypotheses R_n . The following quantities are derived to summarize the results of a MTP: rejection regions (i.e., cut-offs) for the test statistics, confidence regions for the parameters of interest, and adjusted p -values.

As demonstrated in the early article of Pollard and van der Laan (2004), a key feature of our proposed MTPs is the *test statistics null distribution* (rather than data generating null distribution) used to obtain rejection regions, confidence regions, and adjusted p -values. Whether testing single or multiple hypotheses, one needs the (joint) distribution of the test statistics in order to derive a procedure that probabilistically controls Type I errors. In practice, however, the true distribution of the test statistics is unknown and replaced by a null distribution. The choice of a suitable null distribution is crucial, in order to ensure that (finite sample or asymptotic) control of the Type I error rate under the *assumed* null distribution does indeed provide the required control under the *true* distribution. This issue is particularly relevant for large-scale testing problems such as those described above in biomedical and genomic research.

Common approaches use a data generating distribution, such as a permutation distribution, that satisfies the *complete null hypothesis* that all null hypotheses are true. Procedures based on such a *data generating null distribution* typically rely on the *subset pivotality* condition stated in Westfall and Young (1993), p. 42–43, to ensure that control under the data generating null distribution does indeed give the desired control under the true data generating distribution. However, the subset pivotality condition is violated in many important testing problems, since a data generating null distribution may result in a joint distribution for the test statistics that has a different dependence structure than their true distribution. In fact, in most problems, there does not even exist a data generating null distribution that correctly specifies the joint distribution of the test statistics corresponding to the true null hypotheses.

Indeed, subset pivotality fails for two types of testing problems that are highly relevant in biomedical and genomic data analysis: tests concerning correlation coefficients and tests concerning regression parameters. Tests of correlation arise, for example, when seeking to discover sets of co-expressed genes based on microarray expression measures. Tests concerning regression parameters in linear, logistic, survival, and other non-linear models, are commonly used, particularly in medical applications, to identify genes or genomic regions associated with a possibly censored outcome (e.g., survival,

tumor class, response to treatment). While subset pivotality holds for some regression models, such as the simple linear model with independent covariates and error terms, it fails for many models used in practice (e.g., linear regression model of Section 3.1 and logistic regression model of Section 4).

1.2 Outline

The present article, inspired by the early work of Pollard and van der Laan (2004), concerns the choice of a test statistics null distribution in multiple testing. Specifically, it investigates the Type I error and power properties of multiple testing procedures based on our general bootstrap null distribution (Dudoit and van der Laan, 2005; Dudoit et al., 2004b; Pollard and van der Laan, 2004) and various parameter-specific bootstrap null distributions (Westfall and Young, 1993). For the purpose of comparing null distributions, we focus on control of the family-wise error rate (FWER), using the single-step maxT procedure, a common cut-off procedure exploiting the joint distribution of the test statistics. Note, however, that each null distribution could be employed with any other MTP, including our stepwise joint augmentation and empirical Bayes procedures, for controlling generalized tail probability (gTP) error rates, $gTP(q, g) = Pr(g(V_n, R_n) > q)$, for an arbitrary function $g(V_n, R_n)$ of the numbers of false positives V_n and rejected hypotheses R_n (Dudoit and van der Laan, 2005; Dudoit et al., 2004a,b; Pollard and van der Laan, 2004; van der Laan et al., 2005, 2004a,b).

Section 2 provides an overview of our general framework for multiple hypothesis testing and our approach to Type I error control and the choice of a test statistics null distribution. Section 3 describes simulation studies comparing test statistics null distributions in two testing scenarios. The first simulation study considers tests for regression parameters in linear models with dependent covariates and error terms and compares our general non-parametric bootstrap test statistics null distribution (Procedure 2) to a bootstrap null distribution which involves resampling residuals (Westfall and Young (1993), Section 3.4.1, p. 106–109). The second simulation study considers tests for correlation coefficients and compares our general non-parametric bootstrap test statistics null distribution (Procedure 2) to a bootstrap null distribution which involves resampling individual variables independently (Westfall and Young (1993), Section 6.3, p. 194). The simulation studies demonstrate that the choice of null distribution can have a substantial impact on the Type I error and power properties of a given multiple

testing procedure, such as the single-step maxT MTP. Section 4 applies the single-step maxT procedure, based on the general non-parametric bootstrap test statistics null distribution of Procedure 2, to a dataset of microRNA (miRNA) expression measures from cancerous and healthy tissues (Lu et al., 2005). The first testing problem concerns parameters in a (non-linear) logistic regression model relating miRNA expression measures to cancer status, while the second concerns pairwise correlations for miRNA expression measures. We identify 102 (66% of the 155 studied) single miRNAs significantly associated with cancer status, as well as hundreds of pairs of miRNAs with significantly correlated profiles across samples. Finally, Section 5 closes with conclusions and a discussion of ongoing efforts.

2 Methods

2.1 Multiple hypothesis testing framework

The present section introduces a general statistical framework for multiple hypothesis testing and discusses in turn the main ingredients of a multiple testing problem. The reader is referred to Dudoit and van der Laan (2005) and Dudoit et al. (2004b) for details.

2.1.1 Data generating distribution and parameters

Consider a *random sample*, $\mathcal{X}_n \equiv \{X_1, \dots, X_n\}$, of n independent and identically distributed (i.i.d.) random variables from a *data generating distribution* P : $X_i \stackrel{i.i.d.}{\sim} P$, $i = 1, \dots, n$. Suppose that the data generating distribution P is an element of a particular *statistical model* \mathcal{M} , i.e., a set of possibly non-parametric distributions, $P \in \mathcal{M}$. Let P_n denote the corresponding *empirical distribution*, which places probability $1/n$ on each realization of X .

Define *parameters* as arbitrary functions of the data generating distribution P : $\Psi(P) = \psi = (\psi(m) : m = 1, \dots, M)$, where $\psi(m) = \Psi(P)(m) \in \mathbb{R}$.

2.1.2 Null and alternative hypotheses

General submodel hypotheses. In order to cover a broad class of testing problems, define M pairs of null and alternative hypotheses in terms of a collection of *submodels*, $\mathcal{M}(m) \subseteq \mathcal{M}$, $m = 1, \dots, M$, for the data generating distribution P . The M *null hypotheses* are defined as $H_0(m) \equiv \mathbb{I}(P \in \mathcal{M}(m))$

and the corresponding *alternative hypotheses* as $H_1(m) \equiv \mathbf{I}(P \notin \mathcal{M}(m))$. Here, $\mathbf{I}(\cdot)$ is the indicator function, equaling 1 if the condition in parentheses is true and 0 otherwise. Thus, $H_0(m)$ is true, i.e., $H_0(m) = 1$, if the data generating distribution P belongs to submodel $\mathcal{M}(m)$; $H_0(m)$ is false otherwise, i.e., $H_0(m) = 0$.

This general submodel representation covers tests of means, quantiles, correlation coefficients, and regression coefficients in linear and non-linear models (e.g., logistic, survival, time-series, and dose-response models).

Parametric hypotheses. In many testing problems, the submodels concern *parameters*, i.e., each null hypothesis may refer to a single parameter, $\psi(m) = \Psi(P)(m) \in \mathbb{R}$. One distinguishes between two types of testing problems for parametric hypotheses, one-sided and two-sided tests.

$$\begin{array}{ll} \text{One-sided tests} & H_0(m) = \mathbf{I}(\psi(m) \leq \psi_0(m)) \\ \text{vs.} & H_1(m) = \mathbf{I}(\psi(m) > \psi_0(m)), \quad m = 1, \dots, M. \end{array}$$

$$\begin{array}{ll} \text{Two-sided tests} & H_0(m) = \mathbf{I}(\psi(m) = \psi_0(m)) \\ \text{vs.} & H_1(m) = \mathbf{I}(\psi(m) \neq \psi_0(m)), \quad m = 1, \dots, M. \end{array}$$

The hypothesized *null values*, $\psi_0(m)$, are frequently zero. For instance, in microarray data analysis, one may be interested in testing the null hypotheses $H_0(m)$ of no differences in mean gene expression measures between two populations of patients or of no correlation between pairs of gene expression profiles.

Sets of true and false null hypotheses. Let $\mathcal{H}_0 = \mathcal{H}_0(P) \equiv \{m : H_0(m) = 1\} = \{m : P \in \mathcal{M}(m)\}$ be the set of $h_0 \equiv |\mathcal{H}_0|$ *true null hypotheses*, where we note that \mathcal{H}_0 depends on the data generating distribution P . Let $\mathcal{H}_1 = \mathcal{H}_1(P) \equiv \mathcal{H}_0^c(P) = \{m : H_1(m) = 1\} = \{m : P \notin \mathcal{M}(m)\}$ be the set of $h_1 \equiv |\mathcal{H}_1| = M - h_0$ *false null hypotheses*, i.e., *true positives*. The goal of a multiple testing procedure is to accurately estimate the set \mathcal{H}_0 , and thus its complement \mathcal{H}_1 , while probabilistically controlling false positives.

Complete null hypothesis. The *complete null hypothesis*, $H_0^C \equiv \prod_{m=1}^M H_0(m) = \prod_{m=1}^M \mathbf{I}(P \in \mathcal{M}(m)) = \mathbf{I}(P \in \cap_{m=1}^M \mathcal{M}(m))$, is true if and only if all M individual null hypotheses $H_0(m)$ are true, i.e., if and only if the data generating distribution P belongs to the intersection $\cap_{m=1}^M \mathcal{M}(m)$ of the M submodels.

2.1.3 Multiple testing procedure

Test statistics. A *testing procedure* is a *data-driven rule* for deciding which null hypotheses should be rejected, i.e., which $H_0(m)$ should be declared false (zero), so that $P \notin \mathcal{M}(m)$. The decisions to reject or not the null hypotheses are based on an M -vector of *test statistics*, $T_n = (T_n(m) : m = 1, \dots, M)$, that are functions $T_n(m) = T(m; X_1, \dots, X_n)$ of the data, X_1, \dots, X_n . Denote the typically unknown (finite sample) *joint distribution* of the test statistics T_n by $Q_n = Q_n(P)$.

For the test of single-parameter null hypotheses, $H_0(m) = I(\psi(m) \leq \psi_0(m))$ or $H_0(m) = I(\psi(m) = \psi_0(m))$, $m = 1, \dots, M$, consider two main types of test statistics, *difference statistics*,

$$T_n(m) \equiv \text{Estimator} - \text{Null value} = \sqrt{n}(\psi_n(m) - \psi_0(m)), \quad (1)$$

and *t-statistics* (i.e., standardized differences),

$$T_n(m) \equiv \frac{\text{Estimator} - \text{Null value}}{\text{Standard error}} = \sqrt{n} \frac{\psi_n(m) - \psi_0(m)}{\sigma_n(m)}. \quad (2)$$

Here, $\hat{\Psi}(P_n) = \psi_n = (\psi_n(m) : m = 1, \dots, M)$ denotes an *estimator* for the parameter $\Psi(P) = \psi = (\psi(m) : m = 1, \dots, M)$ and $(\sigma_n(m)/\sqrt{n} : m = 1, \dots, M)$ denote the estimated *standard errors* for components $\psi_n(m)$ of ψ_n . Test statistics for other types of null hypotheses include χ^2 -statistics, F -statistics, and likelihood ratio statistics.

Multiple testing procedure. A *multiple testing procedure* (MTP) provides *rejection regions*, $\mathcal{C}_n(m)$, i.e., sets of values for each test statistic $T_n(m)$ that lead to the decision to reject the corresponding null hypothesis $H_0(m)$ and declare that $P \notin \mathcal{M}(m)$, $m = 1, \dots, M$. In other words, a MTP produces a random (i.e., data-dependent) subset \mathcal{R}_n of rejected hypotheses that estimates \mathcal{H}_1 , the set of true positives,

$$\mathcal{R}_n = \mathcal{R}(T_n, Q_{0n}, \alpha) \equiv \{m : H_0(m) \text{ is rejected}\} = \{m : T_n(m) \in \mathcal{C}_n(m)\}, \quad (3)$$

where $\mathcal{C}_n(m) = \mathcal{C}(m; T_n, Q_{0n}, \alpha)$, $m = 1, \dots, M$, denote possibly random rejection regions. The long notation $\mathcal{R}(T_n, Q_{0n}, \alpha)$ and $\mathcal{C}(m; T_n, Q_{0n}, \alpha)$ emphasizes that the MTP depends on:

1. the *data*, $\mathcal{X}_n = \{X_i : i = 1, \dots, n\}$, through the M -vector of *test statistics*, $T_n = (T_n(m) : m = 1, \dots, M)$;

2. an M -variate (estimated) test statistics *null distribution*, Q_{0n} , for deriving rejection regions, confidence regions, and adjusted p -values;
3. the *nominal level* α of the MTP, i.e., the desired upper bound for a suitably defined Type I error rate.

Rejection regions. Rejection regions are typically defined in terms of intervals, such as, $\mathcal{C}_n(m) = (u_n(m), +\infty)$, $\mathcal{C}_n(m) = (-\infty, l_n(m))$, and $\mathcal{C}_n(m) = (-\infty, l_n(m)) \cup (u_n(m), +\infty)$, where $l_n(m) = l(m; T_n, Q_{0n}, \alpha)$ and $u_n(m) = u(m; T_n, Q_{0n}, \alpha)$ are to-be-determined lower and upper *critical values*, or *cut-offs*, computed under the null distribution Q_{0n} for the test statistics T_n . Rejection regions of the form $\mathcal{C}_n(m) = (-\infty, l_n(m)) \cup (u_n(m), +\infty)$ allow the use of asymmetric cut-offs for two-sided tests. Unless specified otherwise, assume that large values of the test statistic $T_n(m)$ provide evidence against the corresponding null hypothesis $H_0(m)$, that is, consider rejection regions of the form $\mathcal{C}_n(m) = (c_n(m), +\infty)$, based on cut-offs $c_n(m) = c(m; T_n, Q_{0n}, \alpha)$. For two-sided tests of single-parameter null hypotheses using difference or t -statistics (Equations (1) and (2)), one could take absolute values of the test statistics.

2.1.4 Type I error rate and power

Type I and Type II errors. In any testing situation, two types of errors can be committed: a *false positive*, or *Type I error*, is committed by rejecting a true null hypothesis, and a *false negative*, or *Type II error*, is committed when the test procedure fails to reject a false null hypothesis. The situation can be summarized by Table 1, below, where the number of rejected hypotheses is $R_n \equiv |\mathcal{R}_n| = \sum_{m=1}^M \mathbf{I}(T_n(m) \in \mathcal{C}_n(m))$, the number of Type I errors is $V_n \equiv |\mathcal{R}_n \cap \mathcal{H}_0| = \sum_{m \in \mathcal{H}_0} \mathbf{I}(T_n(m) \in \mathcal{C}_n(m))$, and the number of Type II errors is $U_n \equiv |\mathcal{R}_n^c \cap \mathcal{H}_1| = \sum_{m \in \mathcal{H}_1} \mathbf{I}(T_n(m) \notin \mathcal{C}_n(m))$. Note that both U_n and V_n depend on the unknown data generating distribution P through the unknown set of true null hypotheses $\mathcal{H}_0 = \mathcal{H}_0(P)$. Therefore, the numbers $h_0 = |\mathcal{H}_0|$ and $h_1 = |\mathcal{H}_1| = M - h_0$ of true and false null hypotheses are *unknown parameters*, the number of rejected hypotheses R_n is an *observable random variable*, and the entries in the body of the table, U_n , $h_1 - U_n$, V_n , and $h_0 - V_n$, are *unobservable random variables*.

Ideally, one would like to simultaneously minimize both the number of Type I errors and the number of Type II errors. Unfortunately, this is not

feasible and one seeks a *trade-off* between the two types of errors. A standard approach is to specify an acceptable level α for a suitably defined Type I error rate and derive testing procedures that aim to minimize a Type II error rate, i.e., maximize power, within the class of tests with Type I error rate at most α .

Type I error rate. When testing multiple hypotheses, there are many possible definitions for the Type I error rate of a test procedure. Accordingly, we adopt a general definition for Type I error rates, as *parameters*, $\theta_n = \theta(F_{V_n, R_n})$, of the joint distribution F_{V_n, R_n} of the numbers of Type I errors $V_n = |\mathcal{R}_n \cap \mathcal{H}_0|$ and rejected hypotheses R_n .

This article focuses on the *family-wise error rate* (FWER), that is, the probability of at least one Type I error,

$$FWER \equiv Pr(V_n > 0) = 1 - F_{V_n}(0). \quad (4)$$

The FWER is controlled, in particular, by the classical Bonferroni procedure (Section 2.3).

Power. As with Type I error rates, *power* can be defined generally as a *parameter*, $\vartheta_n = \vartheta(F_{U_n, R_n})$, of the joint distribution F_{U_n, R_n} of the numbers of Type II errors $U_n = |\mathcal{R}_n^c \cap \mathcal{H}_1|$ and rejected hypotheses R_n .

The present article assesses multiple testing procedures in terms of their *average power*, or expected proportion of rejected false null hypotheses,

$$AvgPwr \equiv \frac{1}{h_1} E[h_1 - U_n] = 1 - \frac{1}{h_1} \int u dF_{U_n}(u). \quad (5)$$

A variety of other Type I and II error rates are discussed in Dudoit and van der Laan (2005).

2.1.5 Adjusted *p*-values

The notion of *p*-value extends directly to multiple testing problems as follows. Consider any multiple testing procedure $\mathcal{R}_n(\alpha) = \mathcal{R}(T_n, Q_{0n}, \alpha)$, with rejection regions $\mathcal{C}_n(m; \alpha) = \mathcal{C}(m; T_n, Q_{0n}, \alpha)$. Then, one can define an *M*-vector of *adjusted p-values*, $\tilde{P}_{0n} = (\tilde{P}_{0n}(m) : m = 1, \dots, M)$, as

$$\begin{aligned} \tilde{P}_{0n}(m) &\equiv \inf \{ \alpha \in [0, 1] : \text{Reject } H_0(m) \text{ at nominal MTP level } \alpha \} \\ &= \inf \{ \alpha \in [0, 1] : m \in \mathcal{R}_n(\alpha) \} \\ &= \inf \{ \alpha \in [0, 1] : T_n(m) \in \mathcal{C}_n(m; \alpha) \}, \quad m = 1, \dots, M. \end{aligned} \quad (6)$$

That is, the adjusted p -value $\tilde{P}_{0n}(m)$, for null hypothesis $H_0(m)$, is the *smallest nominal Type I error level* of the multiple hypothesis testing procedure (e.g., FWER or any other Type I error rate) at which one would reject $H_0(m)$, given T_n . Note that the *unadjusted p -value* $P_{0n}(m)$, for the individual test of null hypothesis $H_0(m)$, corresponds to the special case $M = 1$.

As in single hypothesis tests, the smaller the adjusted p -value, the stronger the evidence against the corresponding null hypothesis. Thus, one rejects $H_0(m)$ for small adjusted p -values $\tilde{P}_{0n}(m)$. This leads to two equivalent representations for a MTP, in terms of rejection regions for the test statistics and in terms of adjusted p -values,

$$\mathcal{R}_n(\alpha) = \{m : T_n(m) \in \mathcal{C}_n(m; \alpha)\} = \{m : \tilde{P}_{0n}(m) \leq \alpha\}. \quad (7)$$

2.2 Type I error rate control and choice of a null distribution

2.2.1 General test statistics null distribution

One of the main tasks in specifying a multiple testing procedure is to derive rejection regions for the test statistics such that the *Type I error rate is controlled* at a desired level α , i.e., such that

$$\begin{aligned} \theta(F_{V_n, R_n}) &\leq \alpha && \text{[finite sample control]} \\ \limsup_{n \rightarrow \infty} \theta(F_{V_n, R_n}) &\leq \alpha && \text{[asymptotic control]}. \end{aligned} \quad (8)$$

Note that the Type I error parameter $\theta(F_{V_n, R_n})$ is defined under the *true distribution* $Q_n = Q_n(P)$ of the test statistics T_n , which is a function of the true underlying data generating distribution P . In practice, however, the distribution $Q_n(P)$ is *unknown and replaced by a null distribution* Q_0 (or estimate thereof, Q_{0n}). The choice of a suitable null distribution Q_0 is crucial, in order to ensure that (finite sample or asymptotic) control of the Type I error rate under this assumed null distribution does indeed provide the required control under the true distribution $Q_n(P)$. For proper control, the null distribution Q_0 must be such that the Type I error rate under this null distribution dominates the Type I error rate under the true distribution $Q_n(P)$. That is, the following *null domination* condition must be satisfied,

$$\begin{aligned} \theta(F_{V_n, R_n}) &\leq \theta(F_{V_0, R_0}) && \text{[finite sample control]} \\ \limsup_{n \rightarrow \infty} \theta(F_{V_n, R_n}) &\leq \theta(F_{V_0, R_0}) && \text{[asymptotic control]}, \end{aligned} \quad (9)$$

where V_0 and R_0 denote, respectively, the numbers of Type I errors and rejected hypotheses under Q_0 , i.e., for $T_n \sim Q_0$.

For error rates $\theta(F_{V_n})$, defined as arbitrary parameters of the distribution of the number of Type I errors V_n , we propose as null distribution $Q_0 = Q_0(P)$, the *asymptotic distribution of the M -vector Z_n of null value shifted and scaled test statistics* (Dudoit and van der Laan, 2005; Dudoit et al., 2004b; Pollard and van der Laan, 2004; van der Laan et al., 2004b),

$$Z_n(m) \equiv \sqrt{\min\left(1, \frac{\tau_0(m)}{\text{Var}[T_n(m)]}\right)} \left(T_n(m) + \lambda_0(m) - E[T_n(m)]\right), \quad m = 1, \dots, M. \quad (10)$$

For the test of single-parameter null hypotheses using t -statistics, the null values are $\lambda_0(m) = 0$ and $\tau_0(m) = 1$. For testing the equality of K population means using F -statistics, the null values are $\lambda_0(m) = 1$ and $\tau_0(m) = 2/(K - 1)$, under the assumption of equal variances in the different populations. Single-step and stepwise procedures based on such a null distribution do indeed provide the desired asymptotic control of the Type I error rate $\theta(F_{V_n})$, for general data generating distributions (with arbitrary dependence structures among variables), null hypotheses (defined in terms of submodels for the data generating distribution), and test statistics (e.g., t -statistics, F -statistics).

For a broad class of testing problems, such as the test of single-parameter null hypotheses using t -statistics (Equation (2)), the null distribution Q_0 is an M -variate Gaussian distribution with mean vector zero and covariance matrix $\Sigma^*(P)$: $Q_0 = Q_0(P) \equiv N(0, \Sigma^*(P))$. For tests of means, where the parameter of interest is the M -dimensional mean vector $\Psi(P) = \psi = E[X]$, the estimator ψ_n is simply the M -vector of empirical means and $\Sigma^*(P)$ is the correlation matrix $\text{Cor}[X]$ of $X \sim P$. More generally, for an asymptotically linear estimator ψ_n , $\Sigma^*(P)$ is the correlation matrix of the vector influence curve (IC). This situation covers standard one-sample and two-sample t -statistics for tests of means, but also test statistics for correlation coefficients (Equation (24)) and regression parameters in linear and non-linear models (Equations (18) and (29)).

In practice, however, since the data generating distribution P is unknown, then so is the proposed null distribution $Q_0 = Q_0(P)$. Resampling procedures, such as the bootstrap procedures of Section 2.4, may be used to conveniently obtain consistent estimators Q_{0n} of the null distribution Q_0 and of the resulting test statistic cut-offs and adjusted p -values (Dudoit and van der

Laan, 2005; Dudoit et al., 2004b; Pollard and van der Laan, 2004; van der Laan et al., 2004b).

2.2.2 Contrast with other approaches

As detailed in Dudoit and van der Laan (2005), Dudoit et al. (2004b), and Pollard and van der Laan (2004), the following two main points distinguish our approach from existing approaches to Type I error rate control and the choice of a null distribution (e.g., in Hochberg and Tamhane (1987) and Westfall and Young (1993)).

Type I error control under the true data generating distribution.

Firstly, we are only concerned with control of the Type I error rate under the *true* data generating distribution P , i.e., under the joint distribution $Q_n = Q_n(P)$ for the test statistics T_n implied by P . The concepts of *weak control* and *strong control* are therefore irrelevant in our context.

In particular, the notion of *null domination* (Dudoit and van der Laan, 2005; Dudoit et al., 2004b) differs from that of *subset pivotality* (Westfall and Young (1993), p. 42–43) in the following senses: (i) null domination is only concerned with the *true* data generating distribution P , i.e., it only considers the subset $\mathcal{H}_0(P)$ of true null hypotheses and not all possible 2^M subsets $\mathcal{J}_0 \subseteq \{1, \dots, M\}$ of null hypotheses, and (ii) null domination does not require equality of the joint distributions Q_{0,\mathcal{H}_0} and $Q_{n,\mathcal{H}_0}(P)$ for the \mathcal{H}_0 -specific test statistics, but the weaker domination of $Q_{n,\mathcal{H}_0}(P)$ by Q_{0,\mathcal{H}_0} .

Null distribution for the test statistics. Secondly, we propose a *null distribution for the test statistics* ($T_n \sim Q_0$) rather than a *data generating null distribution* ($X \sim P_0$). A common choice of data generating null distribution P_0 is one that satisfies the complete null hypothesis, $H_0^C = \mathbf{I}(P \in \cap_{m=1}^M \mathcal{M}(m))$, that all M null hypotheses are true, i.e., $P_0 \in \cap_{m=1}^M \mathcal{M}(m)$. The data generating null distribution P_0 then implies a null distribution $Q_n(P_0)$ for the test statistics.

As discussed in Pollard and van der Laan (2004), procedures based on $Q_n(P_0)$ do not necessarily provide proper (asymptotic) Type I error control under the true distribution P . Indeed, the assumed null distribution $Q_n(P_0)$ and the true distribution $Q_n(P)$ for the test statistics T_n may converge to distributions with different dependence structures and, as a result, may violate the required null domination condition for the Type I error rate (Equation

(9)). For instance, for test statistics with Gaussian asymptotic distributions, the \mathcal{H}_0 -specific correlation matrix under the true distribution P may be different from the corresponding correlation matrix under the assumed complete null distribution P_0 , that is, one may have $\Sigma_{\mathcal{H}_0}^*(P) \neq \Sigma_{\mathcal{H}_0}^*(P_0)$. In the two-sample testing problem, for the commonly-used permutation null distribution P_0 , Pollard and van der Laan (2004) show that $\Sigma_{\mathcal{H}_0}^*(P) = \Sigma_{\mathcal{H}_0}^*(P_0)$ only if (i) the two populations have the same covariance matrices or (ii) the two sample sizes are equal.

Consequently, approaches based on permutation or other data generating null distributions P_0 (e.g., Korn et al. (2004), Troendle (1995, 1996), and Westfall and Young (1993)) are only valid under certain assumptions for the true data generating distribution P . In fact, in most testing problems, there does not exist a data generating null distribution $P_0 \in \cap_{m=1}^M \mathcal{M}(m)$ that correctly specifies a joint distribution for the test statistics, i.e., such that the required null domination condition for the Type I error rate is satisfied.

Thus, unlike current procedures which can only be applied to a limited set of multiple testing problems, the test statistics null distribution Q_0 of Equation (10) leads to single-step and stepwise procedures that provide the desired (asymptotic) Type I error rate control for general data generating distributions, null hypotheses, and test statistics. The null distribution Q_0 can be used in testing problems which cannot be handled by traditional approaches based on a data generating null distribution P_0 and the associated assumption of subset pivotality. Such problems include tests for correlation coefficients and regression parameters in models where covariates and error terms are allowed to be dependent (Sections 3 and 4).

2.3 Multiple testing procedures

The classical *single-step Bonferroni procedure* is perhaps the most widely-used procedure for controlling the family-wise error rate. For a test at nominal FWER level $\alpha \in [0, 1]$, the procedure rejects any hypothesis $H_0(m)$ with unadjusted p -value $P_{0n}(m)$ less than or equal to the common single-step cut-off α/M . The corresponding adjusted p -values are given by,

$$\tilde{P}_{0n}(m) = \min(M P_{0n}(m), 1), \quad m = 1, \dots, M. \quad (11)$$

While simple, this *marginal* procedure can be very conservative for even moderate numbers M of hypotheses. As illustrated in Dudoit et al. (2004a,

2003) and van der Laan et al. (2005), substantial gains in power can be achieved by taking into account the *joint* distribution of the test statistics, as in the following procedure.

Procedure 1 [FWER-controlling single-step maxT procedure] *The single-step maxT procedure is a joint common-cut-off procedure based on the distribution of the maximum test statistic, $\max_{m \in \{1, \dots, M\}} Z(m)$, for an M -vector $Z = (Z(m) : m = 1, \dots, M) \sim Q_0$ with test statistics null distribution Q_0 . For controlling the FWER at nominal level $\alpha \in [0, 1]$, the common cut-off $c(Q_0, \alpha)$, for the test statistics $T_n = (T_n(m) : m = 1, \dots, M)$, is the $(1 - \alpha)$ -quantile of the distribution of $\max_{m \in \{1, \dots, M\}} Z(m)$ under Q_0 ,*

$$c(Q_0, \alpha) \equiv \inf \left\{ c \in \mathbb{R} : \Pr_{Q_0} \left(\max_{m \in \{1, \dots, M\}} Z(m) \leq c \right) \geq (1 - \alpha) \right\}. \quad (12)$$

The corresponding adjusted p -values are given by

$$\tilde{P}_{0n}(m) = \Pr_{Q_0} \left(\max_{m \in \{1, \dots, M\}} Z(m) \geq T_n(m) \right), \quad m = 1, \dots, M. \quad (13)$$

For a test at nominal FWER level α , one has two equivalent representations of the set $\mathcal{R}_n(\alpha)$ of rejected hypotheses, in terms of cut-offs for the test statistics and in terms of adjusted p -values,

$$\mathcal{R}_n(\alpha) = \{m : T_n(m) > c(Q_0, \alpha)\} = \{m : \tilde{P}_{0n}(m) \leq \alpha\}.$$

The reader is referred to our earlier articles and book in preparation, for a variety of other joint multiple testing procedures, controlling a broad class of Type I error rates defined as tail probabilities and expected values for arbitrary functions $g(V_n, R_n)$ of the numbers of false positives V_n and rejected hypotheses R_n (Dudoit and van der Laan, 2005; Dudoit et al., 2004a,b; Pollard and van der Laan, 2004; van der Laan et al., 2005, 2004a,b).

2.4 Bootstrap-based multiple testing procedures

The test statistics null distribution $Q_0 = Q_0(P)$ defined in Equation (10) depends on the true data generating distribution P and is therefore typically unknown. It can be estimated with the (non-parametric or model-based) bootstrap as detailed in Procedure 2, below. Bootstrap-based test statistic cut-offs and adjusted p -values for FWER-controlling single-step maxT Procedure 1 may then be obtained as in Procedure 3.

Procedure 2 [Bootstrap estimation of the test statistics null distribution Q_0] Let P_n^* denote an estimator of the true data generating distribution P . For the non-parametric bootstrap, P_n^* is simply the empirical distribution P_n , that is, samples of size n are drawn at random, with replacement from the observed data, $\mathcal{X}_n = \{X_i : i = 1, \dots, n\}$. For the model-based bootstrap, P_n^* belongs to a model \mathcal{M} for the data generating distribution P , such as a family of multivariate Gaussian distributions. One then proceeds as follows to generate the bootstrap test statistics null distribution.

1. Obtain the b th bootstrap dataset, $\mathcal{X}_n^b \equiv \{X_i^b : i = 1, \dots, n\}$, $b = 1, \dots, B$, by generating n i.i.d. random variables X_i^b with distribution P_n^* .
2. For each bootstrap dataset \mathcal{X}_n^b , compute an M -vector of test statistics, $T_n(\cdot, b) = (T_n(m, b) : m = 1, \dots, M)$, which can be arranged in an $M \times B$ matrix, $\mathbf{T}_n \equiv (T_n(m, b))$, with rows corresponding to the M null hypotheses and columns to the B bootstrap samples.
3. For each null hypothesis $H_0(m)$, compute empirical means $E[T_n(m, \cdot)] \equiv \sum_b T_n(m, b)/B$ and variances $\text{Var}[T_n(m, \cdot)] \equiv \sum_b (T_n(m, b) - E[T_n(m, \cdot)])^2/B$ of the B bootstrap test statistics $T_n(m, b)$ (i.e., row means and variances of the matrix \mathbf{T}_n), to yield estimates of $E[T_n(m)]$ and $\text{Var}[T_n(m)]$, respectively, $m = 1, \dots, M$.
4. Obtain an $M \times B$ matrix, $\mathbf{Z}_n \equiv (Z_n(m, b))$, of null value shifted and scaled bootstrap statistics $Z_n(m, b)$, as in Equation (10), by row-shifting and scaling the matrix \mathbf{T}_n using the bootstrap estimates of $E[T_n(m)]$ and $\text{Var}[T_n(m)]$ and the user-supplied null values $\lambda_0(m)$ and $\tau_0(m)$. That is,

$$Z_n(m, b) \equiv \sqrt{\min \left(1, \frac{\tau_0(m)}{\text{Var}[T_n(m, \cdot)]} \right)} \left(T_n(m, b) + \lambda_0(m) - E[T_n(m, \cdot)] \right). \quad (14)$$

5. The bootstrap estimate Q_{0n} of the null distribution Q_0 from Equation (10) is the empirical distribution of the B columns $Z_n(\cdot, b)$ of matrix \mathbf{Z}_n .

Procedure 3 [Bootstrap estimation of common cut-offs and adjusted p -values for single-step maxT Procedure 1]

0. Apply Procedure 2 to generate an $M \times B$ matrix, $\mathbf{Z}_n = (Z_n(m, b))$, of null value shifted and scaled bootstrap statistics $Z_n(m, b)$.
1. Compute the maximum statistic, $\max_m Z_n(m, b)$, $b = 1, \dots, B$, for each bootstrap dataset \mathcal{X}_n^b , i.e., each column of the matrix \mathbf{Z}_n .
2. For controlling the FWER at nominal level $\alpha \in [0, 1]$, the bootstrap single-step maxT common cut-off $c(Q_{0n}, \alpha)$ is the $(1 - \alpha)$ -quantile of the empirical distribution of the B maxima $\{\max_m Z_n(m, b) : b = 1, \dots, B\}$.
3. The bootstrap single-step maxT adjusted p-value for null hypothesis $H_0(m)$ is the proportion of maxima $\{\max_m Z_n(m, b) : b = 1, \dots, B\}$ exceeding the corresponding observed test statistic $T_n(m)$,

$$\tilde{P}_{0n}(m) \equiv \frac{1}{B} \sum_{b=1}^B \mathbf{I}(\max_m Z_n(m, b) \geq T_n(m)), \quad m = 1, \dots, M. \quad (15)$$

Note that Procedure 3 could be applied, as in Section 3, below, to any matrix \mathbf{Z}_n of resampled statistics (e.g., from other bootstrap or permutation procedures).

3 Simulation studies

This section presents two separate simulation studies comparing our general *non-parametric* bootstrap test statistics null distribution (Procedure 2; Dudoit and van der Laan (2005); Dudoit et al. (2004b); Pollard and van der Laan (2004); van der Laan et al. (2004b)) to *parameter-specific* bootstrap null distributions proposed in Westfall and Young (1993). Specifically, the first simulation study considers tests for *regression parameters* in linear models where the error term is allowed to depend on the covariates and compares the bootstrap null distribution of Procedure 2 to a bootstrap null distribution which involves *resampling residuals* (Westfall and Young (1993), Section 3.4.1, p. 106–109). The second simulation study considers tests for *correlation coefficients* and compares the bootstrap null distribution of Procedure 2 to a bootstrap null distribution which involves *resampling individual variables independently* (Westfall and Young (1993), Section 6.3, p. 194). For both testing problems and each null distribution, the resampling version,

in Procedure 3, of single-step maxT Procedure 1, is applied to control the family-wise error rate.

As detailed in Sections 3.1.4 and 3.2.4, the simulation results demonstrate that the choice of null distribution can have a substantial impact on the Type I error and power properties of a given multiple testing procedure, such as the single-step maxT MTP. The general non-parametric bootstrap test statistics null distribution of Procedure 2 typically controls the Type I error rate “on target” at the nominal level α . In contrast, bootstrapping residuals for the test of regression parameters can lead to severely anti-conservative procedures, while the independent bootstrap for the test of correlation coefficients can lead to conservative procedures.

3.1 Simulation Study 1: Tests of linear regression coefficients in models with dependent covariates and error terms

The first simulation study concerns tests for regression parameters in linear models where the error term is allowed to depend on the covariates. This represents an important and practical testing scenario, since in many biomedical applications, error terms and covariates cannot be assumed to be independent and may have an unknown and complex joint distribution (e.g., logistic model relating cancer status to miRNA expression measures in Section 4).

3.1.1 Simulation model

Data generating distribution. Consider a data structure $(X, Y) \sim P$, where X is an M -dimensional covariate vector and Y a univariate outcome. Assume that the pair (X, Y) has an $(M + 1)$ -dimensional Gaussian distribution P , that satisfies

$$\begin{aligned} E[X] &= 0, & Cov[X] &= \sigma_{xx}, \\ E[Y|X] &= X\psi, & Var[Y|X] &= \sigma_{y|x} = s(X), \end{aligned} \quad (16)$$

where ψ is an M -dimensional vector of regression parameters, σ_{xx} an $M \times M$ covariance matrix, and $s(X)$ a scalar function of the covariates X . That is, one can express the outcome Y in terms of the familiar *linear regression model*

$$Y = X\psi + \epsilon, \quad \text{where} \quad \epsilon|X \sim N(0, s(X)), \quad (17)$$

so that,

$$Y|X \sim N(X\psi, s(X)).$$

Suppose one has a random sample $\mathcal{XY}_n \equiv \{(X_i, Y_i) : i = 1, \dots, n\}$, of n independent and identically distributed pairs $(X_i, Y_i) \sim P$, from the above specified Gaussian data generating distribution P . Let \mathbf{X}_n and \mathbf{Y}_n denote, respectively, the $n \times M$ design matrix and the $n \times 1$ outcome vector.

Null and alternative hypotheses. The hypotheses of interest concern the M components of the regression parameter vector ψ . Specifically, consider two-sided tests of the M null hypotheses $H_0(m) = I(\psi(m) = \psi_0(m))$ vs. the alternative hypotheses $H_1(m) = I(\psi(m) \neq \psi_0(m))$, $m = 1, \dots, M$. For simplicity, set the null values $\psi_0(m)$ equal to zero, i.e., $\psi_0(m) \equiv 0$.

3.1.2 Multiple testing procedures

Test statistics. The M null hypotheses are tested based on standard t -statistics for ordinary least squares (OLS) regression,

$$T_n(m) \equiv \frac{\psi_n(m)}{\sigma_n(m)}, \quad m = 1, \dots, M, \quad (18)$$

where $\psi_n = (\psi_n(m) : m = 1, \dots, M)$ is an M -vector of least squares estimators for the regression parameters, with estimated $M \times M$ covariance matrix σ_n ,

$$\begin{aligned} \psi_n &\equiv (\mathbf{X}_n^\top \mathbf{X}_n)^{-1} \mathbf{X}_n^\top \mathbf{Y}_n, \\ \sigma_n &\equiv \frac{(\mathbf{Y}_n - \mathbf{X}_n \psi_n)^\top (\mathbf{Y}_n - \mathbf{X}_n \psi_n)}{n - M} (\mathbf{X}_n^\top \mathbf{X}_n)^{-1}. \end{aligned} \quad (19)$$

Define an n -vector \mathbf{e}_n of residuals by

$$\mathbf{e}_n \equiv \mathbf{Y}_n - \mathbf{X}_n \psi_n = (e_i \equiv Y_i - X_i \psi_n : i = 1, \dots, n). \quad (20)$$

The simulation study compares the Type I error and power properties of FWER-controlling single-step maxT Procedure 1, based on the following two different bootstrap test statistics null distributions ($B = 10,000$ bootstrap samples).

Bootstrap XY null distribution — Bootstrapping covariate/outcome pairs (X, Y) . The general non-parametric bootstrap test statistics null distribution of Procedure 2 involves *resampling covariate/outcome pairs* (X_i, Y_i) and computing *null value shifted and scaled test statistics* for each bootstrap sample. Specifically, one proceeds as follows for the b th bootstrap sample, $b = 1, \dots, B$.

1. Sample n covariate/outcome pairs (X_i^b, Y_i^b) at random, with replacement from the set of n observations $\mathcal{XY}_n = \{(X_i, Y_i) : i = 1, \dots, n\}$. Let $\mathcal{XY}_n^b \equiv \{(X_i^b, Y_i^b) : i = 1, \dots, n\}$ denote the resulting bootstrap dataset.
2. Compute an M -vector $T_n(\cdot, b) = (T_n(m, b) : m = 1, \dots, M)$ of bootstrap test statistics as in Equation (18), based on the bootstrap dataset \mathcal{XY}_n^b .
3. Compute an M -vector $Z_n(\cdot, b) = (Z_n(m, b) : m = 1, \dots, M)$ of bootstrap null value shifted and scaled test statistics,

$$Z_n(m, b) \equiv \sqrt{\min\left(1, \frac{1}{\text{Var}[T_n(m, \cdot)]}\right)} \left(T_n(m, b) - E[T_n(m, \cdot)]\right),$$

where $E[T_n(m, \cdot)] \equiv \sum_b T_n(m, b)/B$ and $\text{Var}[T_n(m, \cdot)] \equiv \sum_b (T_n(m, b) - E[T_n(m, \cdot)])^2/B$ denote, respectively, the empirical mean and variance of the B bootstrap test statistics $T_n(m, b)$ for null hypothesis $H_0(m)$, $m = 1, \dots, M$ (i.e., row means and variances of the matrix \mathbf{T}_n , as in Procedure 2).

The test statistics null distribution is the empirical distribution Q_{0n} of the $B = 10,000$ M -vectors $\{Z_n(\cdot, b) : b = 1, \dots, B\}$, i.e., of the columns of matrix \mathbf{Z}_n .

Bootstrap e null distribution — Bootstrapping residuals e . In contrast, the parameter-specific bootstrap test statistics null distribution proposed in Section 3.4.1, p. 106–109, of Westfall and Young (1993), involves *resampling residuals* e_i and computing *raw test statistics* (without shifting and scaling) for each bootstrap sample. Specifically, one proceeds as follows for the b th bootstrap sample, $b = 1, \dots, B$.

1. Sample n residuals at random, with replacement from the set of n observed residuals $\{e_i : i = 1, \dots, n\}$ defined in Equation (20). Let $\mathbf{e}_n^b = (e_i^b : i = 1, \dots, n)$ denote the resulting n -vector of bootstrap residuals.
2. Generate n bootstrap covariate/outcome pairs, by randomly pairing each of the n observed covariate vectors X_i with a bootstrap residual e_i^b , that is, by defining a bootstrap outcome n -vector $\mathbf{Y}_n^b \equiv \mathbf{e}_n^b$ as the vector of bootstrap residuals. Let $\mathcal{XY}_n^b \equiv \{(X_i, Y_i^b) : i = 1, \dots, n\}$ denote the resulting bootstrap dataset.
3. Compute an M -vector $T_n(\cdot, b) = (T_n(m, b) : m = 1, \dots, M)$ of bootstrap test statistics as in Equation (18), based on the bootstrap dataset \mathcal{XY}_n^b .

The test statistics null distribution is the empirical distribution Q_{0n} of the $B = 10,000$ M -vectors $\{T_n(\cdot, b) : b = 1, \dots, B\}$, i.e., of the columns of matrix \mathbf{T}_n .

Thus, bootstrap procedures **Bootstrap XY** and **Bootstrap e** differ in two key aspects: (i) the (re)sampling units: **Bootstrap XY** resamples covariate/outcome pairs (X_i, Y_i) , while **Bootstrap e** resamples residuals e_i ; (ii) the bootstrap test statistics: **Bootstrap XY** relies on null value shifted and scaled test statistics Z_n , while **Bootstrap e** relies on “raw” test statistics T_n . In other words, procedure **Bootstrap e** derives the test statistics null distribution by first creating a data generating null distribution in (i), that corresponds to the complete null hypothesis that the outcome Y is independent of each covariate $X(j)$. Note that bootstrapping covariate/outcome pairs (X_i, Y_i) preserves the correlation structure of the data, while bootstrapping residuals and randomly pairing residuals and covariates destroys this correlation.

Single-step maxT procedure. Adjusted p -values for single-step maxT Procedure 1 may be obtained by applying Procedure 3 with bootstrap null distributions **Bootstrap XY** and **Bootstrap e**. Specifically, adjusted p -values for **Bootstrap XY** and **Bootstrap e**, are computed, respectively, from the empirical distributions of the B maxima of shifted and scaled test statistics $\{\max_m Z_n(m, b) : b = 1, \dots, B\}$ and raw test statistics $\{\max_m T_n(m, b) : b = 1, \dots, B\}$. For a test at nominal FWER level α , one rejects null hypotheses with adjusted p -values less than or equal to α .

3.1.3 Simulation study design

Simulation parameters. The following model parameters are varied in the simulation study.

- *Sample size, n .* $n = 25, 100$.
- *Number of hypotheses, M .* $M = 10, 20$.
- *Covariance matrix of the covariates, σ_{xx} .* The covariance matrix σ_{xx} of the covariates has unit diagonal elements and off-diagonal elements set to a common value ς , i.e., $\sigma_{xx}(j, j) = 1$, for $j = 1, \dots, J$, and $\sigma_{xx}(j, j') = \varsigma$, for $j \neq j' = 1, \dots, J$. The following values are considered for the common covariance: $\varsigma = 0.10, 0.50, 0.80$.
- *Conditional variance of outcome Y given covariates X , $s(X)$.* $\text{Var}[Y|X] = \sigma_{y|x} = s(X) = \sum_{m \notin \mathcal{H}_0} X(m)$.
- *Proportion of true null hypotheses, $\frac{h_0}{M}$.* $\frac{h_0}{M} = 0.50, 0.75$.
- *Alternative regression parameters, $(\psi(m) : m \notin \mathcal{H}_0)$.* For each simulation model, regression parameters $(\psi(m) : m \notin \mathcal{H}_0)$, for the false null hypotheses, are generated as $|\mathcal{H}_0^c| = M - h_0$ independent uniform random variables over the interval $[0, \frac{\mu}{\sqrt{n}}]$. That is, $\psi(m) \stackrel{i.i.d.}{\sim} U(0, \frac{\mu}{\sqrt{n}})$, $m \notin \mathcal{H}_0$. The following values are considered for the magnitude parameter: $\mu = 0.10, 0.25$.

Estimating Type I error rate and power. For each simulation model (i.e., each combination of parameter values n , M , ς , $s(X)$, h_0/M , and μ), generate $A = 500$ random samples, $\mathcal{XY}_n^a \equiv \{(X_i^a, Y_i^a) : i = 1, \dots, n\}$, of covariate/outcome pairs $(X, Y) \sim P$. For each such simulated dataset, compute adjusted p -values $\tilde{P}_{0n}^a(m)$ for single-step maxT Procedure 3, based on each of the two bootstrap null distributions (**Bootstrap XY** and **Bootstrap e**). For a given nominal Type I error level α , compute the numbers of rejected

hypotheses $R_n^a(\alpha)$, Type I errors $V_n^a(\alpha)$, and Type II errors $U_n^a(\alpha)$,

$$\begin{aligned} R_n^a(\alpha) &\equiv \sum_{m=1}^M \mathbf{I}(\tilde{P}_{0n}^a(m) \leq \alpha), \\ V_n^a(\alpha) &\equiv \sum_{m \in \mathcal{H}_0} \mathbf{I}(\tilde{P}_{0n}^a(m) \leq \alpha), \\ U_n^a(\alpha) &\equiv \sum_{m \notin \mathcal{H}_0} \mathbf{I}(\tilde{P}_{0n}^a(m) > \alpha). \end{aligned} \quad (21)$$

The *actual Type I error rate* is estimated as follows and then compared to the *nominal Type I error level* α ,

$$FWER(\alpha) \equiv \frac{1}{A} \sum_{a=1}^A \mathbf{I}(V_n^a(\alpha) > 0). \quad (22)$$

The *average power* of a given MTP is estimated by

$$AvgPwr(\alpha) \equiv 1 - \frac{1}{h_1} \frac{1}{A} \sum_{a=1}^A U_n^a(\alpha). \quad (23)$$

The simulation error for the actual Type I error rate and power is of the order $1/\sqrt{A} = 1/\sqrt{500} \approx 0.045$.

Graphical summaries. Simulation results are displayed using the following two main types of graphical summaries.

- **Type I error control comparison.** For a given data generating model, plot, for each MTP, the *difference between the nominal and actual Type I error rates* vs. the *nominal Type I error rate*, i.e., plot

$$(\alpha - FWER(\alpha)) \quad \text{vs.} \quad \alpha,$$

for $\alpha \in \{0, 0.01, 0.02, \dots, 0.50\}$, i.e., values of α in `seq(from = 0, to = 0.50, by=0.01)`. Positive (negative) differences correspond to (anti-)conservative MTPs; the higher the curve, the more conservative the procedure.

- **Power comparison.** For a given data generating model, *receiver operator characteristic* (ROC) curves, comparing different MTPs in terms

of power, can be obtained by plotting, for each MTP, *power* vs. *actual Type I error rate*, i.e., $AvgPwr(\alpha)$ vs. $FWER(\alpha)$, for a range of nominal Type I error levels α . However, due to possibly large variations in power between simulation models, consider instead the following modified display, which facilitates comparisons across models. For a given model, plot the *difference in power* between two procedures vs. the *actual Type I error rate*, i.e., plot

$$(AvgPwr^{Boot\ XY}(\alpha) - AvgPwr^{Boot\ e}(\alpha)) \quad \text{vs.} \quad FWER(\alpha),$$

for $\alpha \in \{0, 0.01, 0.02, \dots, 0.50\}$.

3.1.4 Simulation results

Our comparison of the test statistics null distributions **Bootstrap XY** and **Bootstrap e** focusses on Type I error control.

Figure 1 displays differences between nominal and actual Type I error rates for four simulation models, where one parameter is varied as the others remain constant. In general, procedures based on the residual bootstrap null distribution **Bootstrap e** are *anti-conservative* over the entire range of the nominal level α , while procedures based on the general non-parametric bootstrap null distribution **Bootstrap XY** control the Type I error rate at the target nominal level α . In some testing scenarios, the actual Type I error rate for **Bootstrap e** exceeds the nominal Type I error level by as much as 0.20. We comment on a number of trends below.

- *Covariance matrix of the covariates, σ_{xx} .* (Figure 1, Panels (b) vs. (a)) As the correlation ς between covariates increases, the actual Type I error rate for **Bootstrap XY** gets closer to the nominal level α . In contrast, procedure **Bootstrap e** becomes more anti-conservative as the correlation ς increases.
- *Sample size, n .* (Figure 1, Panels (c) vs. (a)) As the number of observations n increases, the actual Type I error rate for **Bootstrap XY** gets closer to the nominal level α .
- *Alternative regression parameters, $(\psi(m) : m \notin \mathcal{H}_0)$.* (Figure 1, Panels (d) vs. (a)) As the magnitude of the parameter μ , defining the regression coefficients $(\psi(m) : m \notin \mathcal{H}_0)$ for the false null hypotheses,

increases, the actual Type I error rate for **Bootstrap XY** gets closer to the nominal level α . In contrast, procedure **Bootstrap e** becomes more anti-conservative as the magnitude μ increases.

- *Proportion of true null hypotheses, $\frac{h_0}{M}$.* No clear trends are noticeable for the proportion of true null hypotheses.

We find that, for most simulation models, the differences in power are within simulation error (i.e., less than $1/\sqrt{A} = 1/\sqrt{500} \approx 0.045$), for the two versions of bootstrap-based single-step maxT Procedure 3 (Figure 2). Procedure **Bootstrap e** tends to be slightly more powerful than **Bootstrap XY**, especially for the larger correlation value $\varsigma = 0.80$. The main noticeable trends are, as expected, that power increases with sample size n and effect size μ .

3.2 Simulation Study 2: Tests of correlation coefficients

The second simulation study concerns tests of correlation coefficients, a testing scenario of great interest in genomic applications. Indeed, as illustrated in Section 4, below, a common question is the identification of genes with correlated expression profiles in microarray and other high-throughput gene expression assays.

3.2.1 Simulation model

Data generating distribution. Consider a random J -vector $X \sim P = N(0, \sigma)$, with J -dimensional Gaussian distribution P , mean vector zero, and covariance matrix $\sigma = (\sigma(j, j') : j, j' = 1, \dots, J)$ equal to the corresponding correlation matrix $\rho = (\rho(j, j') : j, j' = 1, \dots, J)$.

Suppose one has a random sample $\mathcal{X}_n \equiv \{X_i : i = 1, \dots, n\}$, of n i.i.d. random variables $X_i \sim P$, from the above specified Gaussian data generating distribution P .

Null and alternative hypotheses. The hypotheses of interest concern the $M \equiv \binom{J}{2} = J(J-1)/2$ distinct entries $\psi = (\psi(m) : m = 1, \dots, M)$ of the $J \times J$ correlation matrix ρ . One may recode pairs of row and column indices $\{(j, j') : j = 1, \dots, (J-1), j' = j+1, \dots, J\}$, for the upper triangle of ρ , into a single index $m = 1, \dots, M$, defined by $m \equiv (j-1)(2J-j)/2 + (j' - j)$.

Consider two-sided tests of the $M = J(J-1)/2$ null hypotheses $H_0(m) = I(\psi(m) = \psi_0(m))$ vs. the alternative hypotheses $H_1(m) = I(\psi(m) \neq \psi_0(m))$, $m = 1, \dots, M$. For simplicity, set the null values $\psi_0(m)$ equal to zero, i.e., test the null hypotheses of no correlation.

3.2.2 Multiple testing procedures

Test statistics. The M null hypotheses are tested based on the following t -statistics,

$$T_n(m) \equiv \sqrt{n-2} \frac{\psi_n(m)}{\sqrt{1 - \psi_n^2(m)}}, \quad m = 1, \dots, M, \quad (24)$$

where $\psi_n = (\psi_n(m) : m = 1, \dots, M)$ is the M -vector of empirical correlation coefficients. Specifically, the empirical correlation coefficient for the pair of random variables $(X(j), X(j'))$, corresponding to the m th null hypothesis, is defined as

$$\psi_n(m) = \rho_n(j, j') \equiv \frac{\sigma_n(j, j')}{\sqrt{\sigma_n(j, j) \sigma_n(j', j')}}, \quad (25)$$

based on empirical means $\bar{X}_n(j)$ and covariances $\sigma_n(j, j')$,

$$\bar{X}_n(j) \equiv \frac{1}{n} \sum_{i=1}^n X_i(j), \quad \sigma_n(j, j') \equiv \frac{1}{n} \sum_{i=1}^n (X_i(j) - \bar{X}_n(j))(X_i(j') - \bar{X}_n(j')).$$

For Gaussian data generating distributions, the t -statistics in Equation (24) have marginal t -distributions with $(n-2)$ degrees of freedom, under the null hypotheses that the corresponding correlation coefficients are zero, i.e., $\psi(m) = 0$ (with the finite sample bias correction $n/(n-1)$ in the definition of the empirical covariance matrix σ_n).

One could also use unstandardized test statistics,

$$T_n(m) \equiv \sqrt{n} \psi_n(m), \quad m = 1, \dots, M. \quad (26)$$

The simulation study compares the Type I error and power properties of FWER-controlling single-step maxT Procedure 1, based on the following two different bootstrap test statistics null distributions ($B = 10,000$ bootstrap samples).

Bootstrap X null distribution — Bootstrapping entire J-vectors X.

The general non-parametric bootstrap test statistics null distribution of Procedure 2 involves *resampling entire J-vectors* X_i and computing *null value shifted and scaled test statistics* for each bootstrap sample. Specifically, one proceeds as follows for the b th bootstrap sample, $b = 1, \dots, B$.

1. Sample n J -vectors X_i^b at random, with replacement from the set of n observations $\mathcal{X}_n = \{X_i : i = 1, \dots, n\}$. Let $\mathcal{X}_n^b \equiv \{X_i^b : i = 1, \dots, n\}$ denote the resulting bootstrap dataset.
2. Compute an M -vector $T_n(\cdot, b) = (T_n(m, b) : m = 1, \dots, M)$ of bootstrap test statistics as in Equation (24), based on the bootstrap dataset \mathcal{X}_n^b .
3. Compute an M -vector $Z_n(\cdot, b) = (Z_n(m, b) : m = 1, \dots, M)$ of bootstrap null value shifted and scaled test statistics,

$$Z_n(m, b) \equiv \sqrt{\min\left(1, \frac{1}{\text{Var}[T_n(m, \cdot)]}\right)} \left(T_n(m, b) - E[T_n(m, \cdot)]\right),$$

where $E[T_n(m, \cdot)] \equiv \sum_b T_n(m, b)/B$ and $\text{Var}[T_n(m, \cdot)] \equiv \sum_b (T_n(m, b) - E[T_n(m, \cdot)])^2/B$ denote, respectively, the empirical mean and variance of the B bootstrap test statistics $T_n(m, b)$ for null hypothesis $H_0(m)$, $m = 1, \dots, M$ (i.e., row means and variances of the matrix \mathbf{T}_n , as in Procedure 2).

The test statistics null distribution is the empirical distribution Q_{0n} of the $B = 10,000$ M -vectors $\{Z_n(\cdot, b) : b = 1, \dots, B\}$, i.e., of the columns of matrix \mathbf{Z}_n .

Bootstrap X(j) null distribution — Bootstrapping independent entries $X(j)$ of the J -vectors X . In contrast, the parameter-specific bootstrap test statistics null distribution proposed in Section 6.3, p. 194, of Westfall and Young (1993), involves *resampling each component* $X_i(j)$ of the J -vectors X_i *independently* and computing *raw test statistics* (without shifting and scaling) for each bootstrap sample. Specifically, one proceeds as follows for the b th bootstrap sample, $b = 1, \dots, B$.

1. For each variable $X(j)$, $j = 1, \dots, J$, sample n j -specific entries $X_i^b(j)$, $i = 1, \dots, n$, at random, with replacement from the set of n j -specific

observations $\{X_i(j) : i = 1, \dots, n\}$. The i th bootstrap J -vector $X_i^b = (X_i^b(j) : j = 1, \dots, J)$, $i = 1, \dots, n$, is obtained by combining J such independently sampled variables. Let $\mathcal{X}_n^b \equiv \{X_i^b : i = 1, \dots, n\}$ denote the resulting bootstrap dataset.

2. Compute an M -vector $T_n(\cdot, b) = (T_n(m, b) : m = 1, \dots, M)$ of bootstrap test statistics as in Equation (24), based on the bootstrap dataset \mathcal{X}_n^b .

The test statistics null distribution is the empirical distribution Q_{0n} of the $B = 10,000$ M -vectors $\{T_n(\cdot, b) : b = 1, \dots, B\}$, i.e., of the columns of matrix \mathbf{T}_n .

As in the regression example of Section 3.1, bootstrap procedures **Bootstrap X** and **Bootstrap X(j)** differ in two key aspects: (i) the (re)sampling units: **Bootstrap X** resamples entire J -vectors X_i , while **Bootstrap X(j)** resamples independent components $X_i(j)$; (ii) the bootstrap test statistics: **Bootstrap X** relies on null value shifted and scaled test statistics Z_n , while **Bootstrap X(j)** relies on “raw” test statistics T_n . In other words, procedure **Bootstrap X(j)** derives the test statistics null distribution by first creating a data generating null distribution in (i), that corresponds to the complete null hypothesis that the J variables $X(j)$, $j = 1, \dots, J$, are independent.

Single-step maxT procedure. Adjusted p -values for single-step maxT Procedure 1 may be obtained by applying Procedure 3 with bootstrap null distributions **Bootstrap X** and **Bootstrap X(j)**. Specifically, adjusted p -values for **Bootstrap X** and **Bootstrap X(j)**, are computed, respectively, from the empirical distributions of the B maxima of shifted and scaled test statistics $\{\max_m Z_n(m, b) : b = 1, \dots, B\}$ and raw test statistics $\{\max_m T_n(m, b) : b = 1, \dots, B\}$. For a test at nominal FWER level α , one rejects null hypotheses with adjusted p -values less than or equal to α .

3.2.3 Simulation study design

Simulation parameters. The following model parameters are used in the simulation study.

- *Sample size, n .* $n = 25$.
- *Number of hypotheses, M .* $M = 45$.

- *Proportion of true null hypotheses*, $\frac{h_0}{M}$. $\frac{h_0}{M} = \frac{25}{45} \approx 0.56$.
- *Correlation matrix*, ρ . The correlation matrix $\rho = (\rho(j, j') : j, j' = 1, \dots, J)$ (here, equal to the covariance matrix σ) has the following block diagonal form,

$$\rho = \begin{bmatrix} \varrho_{J/2 \times J/2} & O_{J/2 \times J/2} \\ O_{J/2 \times J/2} & \varrho_{J/2 \times J/2} \end{bmatrix}, \quad (27)$$

where $O_{J/2 \times J/2}$ denotes a $J/2 \times J/2$ matrix of zeros and $\varrho_{J/2 \times J/2}$ a $J/2 \times J/2$ matrix with unit diagonal elements and off-diagonal elements set to a common value ϱ , i.e., $\varrho_{J/2 \times J/2}(j, j) = 1$, for $j = 1, \dots, J/2$, and $\varrho_{J/2 \times J/2}(j, j') = \varrho$, for $j \neq j' = 1, \dots, J/2$. The following values are considered for the common block correlation coefficient: $\varrho = 0.30, 0.50, 0.60$.

Note that the only parameter that is varied in the simulation study is the correlation matrix, ρ , that is, the parameter of interest in the multiple testing problem.

Estimating Type I error rate and power. As in Section 3.1.3, above, for Simulation Study 1.

Graphical summaries. As in Section 3.1.3, above, for Simulation Study 1.

3.2.4 Simulation results

Our comparison of the test statistics null distributions **Bootstrap X** and **Bootstrap X(j)** focusses on Type I error rate control.

Figure 3 displays differences between nominal and actual Type I error rates for three simulation models, where the common block correlation coefficient ϱ is varied as the other parameters remain constant. In general, procedures based on the independent covariates bootstrap null distribution **Bootstrap X(j)** are *conservative* over the entire range of the nominal level α , while procedures based on the general non-parametric bootstrap null distribution **Bootstrap X** control the Type I error rate closer to the target nominal level α . The most extreme differences are observed for large nominal Type I error levels α . In some testing scenarios, the nominal level for **Bootstrap X(j)**

exceeds the actual Type I error rate by as much as 0.25. As the correlation parameter ρ increases, procedure **Bootstrap X(j)** becomes more conservative.

As in the first simulation study, we find that, for most simulation models, the differences in power are within simulation error (i.e., less than $1/\sqrt{A} = 1/\sqrt{500} \approx 0.045$), for the two versions of bootstrap-based single-step maxT Procedure 3 (Figure 4). Procedure **Bootstrap X** tends to be slightly more powerful than **Bootstrap X(j)**, especially for larger correlation parameters ρ . The main noticeable trends are, as expected, that power increases with sample size n and effect size ρ .

Similar trends are observed for standardized (Equation (24)) and unstandardized (Equation (26)) correlation test statistics (data not shown for unstandardized statistics).

4 microRNA data analysis

In addition to playing the important role of passing genetic messages from DNA to the protein-making machinery of the cell, RNA serves many other cellular functions. A new class of small, non-coding RNAs, known as *microRNAs* (miRNA), are currently the subject of intense study due to their provocative roles in controlling developmental timing and regulating the translation of *messenger RNA* (mRNA). Like transcription factor proteins, miRNAs have the potential to affect the abundance of a wide range of proteins in the cell. In their recent article, investigating miRNA levels in cells derived from cancerous and non-cancerous tissues, Lu et al. (2005) made an astonishing discovery: predictors based on abundance of the several hundred known mammalian miRNAs are better able to distinguish developmental lineage, differentiation state, and cancer state, than the best predictors based on genome-wide mRNA expression profiles from the same cells. Motivated by these findings, we have undertaken further analysis of this publicly available miRNA dataset.

4.1 miRNA dataset of Lu et al. (2005)

Lu et al. (2005) measured expression levels of 217 known human miRNAs, by a bead-based flow cytometric profiling method, in cells from 46 cancerous and 140 healthy tissues ($n = 186$ total samples). The pre-processed, \log_2 -transformed data are available from the authors' website (www.broad.mit.edu).

edu/cancer/pub/miGCM: miRNA expression measures in file `miGCM_218.gct`; probe sequence information in file `supplementary_table_1.xls`; target sample information, such as cancer status and tissue type, in file `supplementary_table_2.xls`). The analyses in Sections 4.2 and 4.3, below, exclude cell lines and any miRNA with expression measures below a detection threshold of $\log_2(32) = 5$ in more than half of the $n = 186$ samples.

The data for each of the $n = 186$ samples consist of a binary outcome Y_i for cancer status (1 for cancerous vs. 0 for healthy tissues) and a $J = 155$ -dimensional covariate vector $X_i = (X_i(j) : j = 1, \dots, J)$ of real-valued expression measures for each of $J = 155$ miRNAs, $i = 1, \dots, n$.

4.2 Tests of logistic regression coefficients: Differentially expressed miRNAs between cancerous and healthy tissues

The original publication of Lu et al. (2005) includes a comparison of miRNA expression measures between cancerous and healthy tissues, using the FWER-controlling marginal Bonferroni procedure, with modified two-sample t -statistics. For a test at nominal FWER level $\alpha = 0.05$, the authors found that 59% of the miRNAs were significantly less abundant in cancerous compared to healthy tissues. Only a few miRNAs were over-expressed in cancerous tissues and none significantly so. miRNA measures were observed to vary greatly among the 19 different tissue types represented in the dataset (stomach, colon, pancreas, etc.); tissue type was therefore treated as a *confounding variable*.

Multiple testing procedure. An alternative approach to identifying *differentially expressed* miRNAs between cancerous and healthy tissues is to fit, for each miRNA, a *logistic regression model* including tissue type as a covariate. Specifically, the logistic regression model for the j th miRNA is

$$\text{logit}(E[Y|X]) \equiv \alpha(j) + \beta(j)X(j) + \gamma(j)W, \quad j = 1, \dots, J, \quad (28)$$

where $\text{logit}(z) \equiv \log(z/(1-z))$ is the *logit function*, W is a 19-dimensional tissue type indicator vector, and $\gamma(j)$ a corresponding miRNA-specific 19-dimensional parameter vector.

The parameter of interest in the logistic model of Equation (28) is $\beta(j)$, the scalar coefficient for the expression measure $X(j)$ of the j th miRNA, $j = 1, \dots, J$. Thus, for each miRNA, one considers two-sided tests of the

null hypotheses $H_0(j) = I(\beta(j) = 0)$, of no association of the expression measures $X(j)$ with cancer status Y , vs. the alternative hypotheses $H_1(j) = I(\beta(j) \neq 0)$. Two-sided tests are used to identify both over- and under-expressed miRNAs in cancerous tissues.

The J null hypotheses are tested based on the following t -statistics,

$$T_n(j) \equiv \frac{\beta_n(j) - \beta_0(j)}{\sigma_n(j)}, \quad j = 1, \dots, J, \quad (29)$$

where the null values $\beta_0(j)$ are zero and $\beta_n(j)$ are logistic regression parameter estimates with estimated standard errors $\sigma_n(j)$ (as implemented in the R function `glm` for the binomial family, with the call `glm(Y ~ X(j) + W, family="binomial")`).

In order to simultaneously test the J null hypotheses of no association of miRNA measures with cancer status, we apply FWER-controlling single-step maxT Procedure 1, with the general non-parametric bootstrap test statistics null distribution of Procedure 2 ($B = 5,000$ bootstrap samples). That is, test statistic cut-offs and adjusted p -values are computed as in Procedure 3. Note that fitting the logistic regression model of Equation (28) allows the identification of differentially expressed miRNAs, while adjusting for tissue type.

Results. Bootstrap-based single-step maxT Procedure 3 yields 102 miRNAs (66% of the 155 studied) with adjusted p -values less than a nominal FWER level $\alpha = 0.05$. In fact, 53 of the p -values are less than 0.01 and 36 are approximately equal to zero, indicating that some miRNAs are very significantly differentially expressed between cancerous and healthy tissues (Table 2 and Figure 5). All of the 102 significant miRNAs have test statistics $T_n(j) < -3.6$, suggesting under-expression in cancerous compared to healthy tissues. These findings are in agreement with the original publication of Lu et al. (2005), but a larger proportion of the studied miRNAs are found to be significantly differentially expressed with our proposed bootstrap method.

Five of the highly significant miRNAs listed in Table 2 are located in minimal deleted regions, minimal amplified regions, and breakpoint regions involved in human cancers (Calin et al., 2004). Specifically, miR-23b and let-7d have been associated with urothelial cancer, miR-100 with breast, lung, ovarian, and cervical cancers, miR-22 with hepatocellular cancer, and miR99a with lung cancer.

Note that another approach for comparing mean miRNA expression measures in cancerous vs. healthy tissues could be based on standard two-sample

t -statistics. For such simple tests, data generating null distributions, such as the permutation distribution, lead to proper Type I error control under the conditions that (i) the two populations have the same covariance matrices or (ii) the two sample sizes are equal (Pollard and van der Laan, 2004). Our general multiple testing methodology, however, allows one to use more general and flexible models, such as the logistic regression model of Equation (28), which facilitates adjustment for covariates and also provides a simple predictor of cancer status.

4.3 Tests of correlation coefficients: Co-expressed miRNAs

A biological question of great interest in gene expression experiments is the identification of *co-expressed* genes, here, miRNAs with correlated expression measures across tissue samples. While some tests of association between miRNA abundance and a binary outcome, such as cancer status, could be performed with standard multiple testing tools (e.g., MTPs based on a permutation data generating null distribution), correlation tests are a case where our bootstrap-based MTPs truly allow one to perform previously unavailable analyses.

Multiple testing procedures. Consider, the $M \equiv J(J - 1)/2 = 155 \times 154/2 = 11,935$ distinct Pearson correlation coefficients between pairs of miRNA expression profiles,

$$\rho(j, k) \equiv \text{Cor}[X(j), X(k)], \quad j = 1, \dots, J - 1, \quad k = j + 1, \dots, J. \quad (30)$$

It is of interest to identify all pairs of miRNAs with significantly correlated expression profiles across the $n = 186$ samples. Thus, for each distinct pair (j, k) of miRNAs, one considers two-sided tests of the null hypotheses $H_0(j, k) = \text{I}(\rho(j, k) = 0)$, of no association in expression measures, vs. the alternative hypotheses $H_1(j, k) = \text{I}(\rho(j, k) \neq 0)$.

The M null hypotheses are tested based on the following test statistics,

$$T_n(j, k) \equiv \sqrt{n} \rho_n(j, k), \quad j = 1, \dots, J - 1, \quad k = j + 1, \dots, J, \quad (31)$$

where $\rho_n(j, k)$ are empirical correlation coefficients, as defined previously in Section 3.2.2.

In order to simultaneously test the $M = 11,935$ null hypotheses of no association in expression measures between pairs of miRNAs, we again apply FWER-controlling single-step maxT Procedure 1, with the general non-parametric bootstrap test statistics null distribution of Procedure 2 ($B = 5,000$ bootstrap samples). That is, test statistic cut-offs and adjusted p -values are computed as in Procedure 3.

Results. Interestingly, bootstrap-based single-step maxT Procedure 3 yields 8,916 miRNA pairs (or nearly 75% of all $M = 11,935$ pairs) with adjusted p -values less than a nominal FWER level $\alpha = 0.05$ and 7,479 with p -values approximately equal to zero (Figure 6). Correlations found to be statistically significantly different from zero at nominal level $\alpha = 0.05$ range from 0.26 to 0.99, with median value 0.55. The most significant are given in Table 3. Several pairs are composed of miRNAs in the same family, e.g., **hsa-miR-10a** and **hsa-miR-10b**. Only 8% of all pairwise correlations are negative, and none are significantly so.

The two most significantly correlated miRNAs are a pair of paralogs, **miR-17-5p** (chromosome 17) and **miR-106a** (chromosome X), which are part of miRNA clusters believed to be up-regulated by the proto-oncogene c-MYC (O'Donnell et al., 2005). **miR-19a**, **miR-19b**, and **miR-20** are also members of these paralogous miRNA clusters. Several other co-expressed miRNAs are linked to cancer. In particular, **miR-107** has been shown to increase cell growth in lung carcinomas (Cheng et al., 2005). **miR-143** and **miR-145**, located within 1.7 kb on human chromosome 5, are expressed at lower levels in cancerous and pre-cancerous tissue compared to normal colon tissue (Michael et al., 2003).

Cluster analysis. The above multiple testing analysis clearly suggests the existence of clusters of highly correlated miRNAs. We therefore decided to perform hierarchical clustering of the miRNAs, in order to identify general expression patterns and groups of co-expressed miRNAs. We use the *hierarchical ordered partitioning and collapsing hybrid* (HOPACH) algorithm with Pearson correlation distance (van der Laan and Pollard, 2003). Figure 7 displays the 155×155 miRNA correlation matrix, ordered according to the final level of the HOPACH tree, so that similarly expressed miRNAs appear near each other. It will be of great interest to investigate the biological and medical implications of clusters of co-expressed miRNAs.

5 Conclusions

This investigation of multiple testing procedures has focused on the choice of a test statistics null distribution, in testing problems for which subset pivotality does not hold. Subset pivotality is a condition under which data generating distributions satisfying the complete null hypothesis produce valid test statistics null distributions (Westfall and Young (1993), p. 42–43). Commonly-used permutation or parameter-specific bootstrap test statistics null distributions rely on the subset pivotality condition to justify Type I error control under the true distribution. However, subset pivotality is violated in many important testing problems, since a data generating null distribution may result in a joint distribution for the test statistics that has a different dependence structure than their true distribution. In fact, in most situations, there does not even exist a data generating null distribution that correctly specifies the joint distribution of the test statistics corresponding to the true null hypotheses.

Indeed, subset pivotality fails for two types of testing problems that are highly relevant in biomedical and genomic data analysis: tests concerning correlation coefficients and tests concerning regression parameters. Correlation tests abound in molecular biology, where similarities between measurable properties of large numbers of genes and genome sequences are of great interest. Non-linear regression models are also frequently used to assess genotype/phenotype associations, while adjusting for potential confounding variables. Procedures based on a data generating null distribution, such as a permutation distribution, do not provide a correct test statistics null distribution in these settings.

Motivated by limitations of existing approaches, Pollard and van der Laan (2004), and subsequently Dudoit et al. (2004b), propose a general characterization and explicit construction of a test statistics null distribution that controls Type I errors, without requirements such as subset pivotality, in testing problems involving general data generating distributions (i.e., arbitrary dependence structures among variables). Resampling procedures, such as the bootstrap procedures of Section 2.4, are provided to conveniently obtain consistent estimators of the null distribution and of the resulting test statistic cut-offs and adjusted p -values. Pollard and van der Laan (2004) compare MTPs based on the proposed bootstrap test statistics null distribution and several other null distributions in the two-sample testing problem. The former null distribution performs competitively whenever the sample

sizes are large enough to avoid ties in the resampled distribution and poorly estimated variances in the denominators of t -statistics.

The goal of the present paper was to evaluate the practical performance of different test statistics null distributions in cases where subset pivotality fails. Specifically, the simulation studies of Section 3 compare our general non-parametric bootstrap test statistics null distribution (Procedure 2) to parameter-specific bootstrap null distributions, in the following two settings: tests for regression coefficients in linear models where covariates and error terms are allowed to be dependent and tests for correlation coefficients. The general non-parametric bootstrap distribution (Procedure 2, **Bootstrap XY** and **Bootstrap X**) differs from corresponding parameter-specific bootstrap distributions (**Bootstrap e** and **Bootstrap X(j)**) in two key aspects: (i) the (re)sampling units: **Bootstrap XY** and **Bootstrap X** resample “raw” observations, while **Bootstrap e** and **Bootstrap X(j)** resample, respectively, residuals e_i and independent components $X_i(j)$; (ii) the bootstrap test statistics: **Bootstrap XY** and **Bootstrap X** rely on null value shifted and scaled test statistics Z_n , while **Bootstrap e** and **Bootstrap X(j)** rely on “raw” test statistics T_n . In other words, procedures **Bootstrap e** and **Bootstrap X(j)** derive the test statistics null distribution by first creating a data generating distribution that satisfies the complete null hypothesis.

The simulation studies, involving a range of data generating models, demonstrate that the choice of null distribution can have a substantial impact on the Type I error and power properties of a given multiple testing procedure. The single-step maxT procedure, based on the general non-parametric bootstrap null distribution of Procedure 2, does indeed control the family-wise error rate at or slightly below the target nominal level. Interestingly, comparable MTPs based on parameter-specific bootstrap null distributions, are anti-conservative for tests of regression coefficients (**Bootstrap e**) and conservative for tests of correlation coefficients (**Bootstrap X(j)**). Power is similar for the different null distributions in both testing problems.

Section 4 illustrates the flexibility and power of our proposed methodology, by applying the single-step maxT procedure, with general non-parametric bootstrap test statistics null distribution (Procedures 2 and 3), to a dataset of miRNA expression measures from cancerous and healthy tissues (Lu et al., 2005). Tests for regression coefficients, in a logistic model adjusting for tissue type, identify 102 miRNAs as being significantly differentially expressed between cancerous and healthy tissues (nominal FWER level 0.05). This corroborates the original article’s discovery that miRNA expression profiling has

great potential for cancer diagnosis. Stepwise, augmentation, and empirical Bayes procedures could be used for more powerful analyses and control of a broader class of Type I error rates (Dudoit and van der Laan, 2005; Dudoit et al., 2004a; van der Laan et al., 2005, 2004a,b).

We also investigated several questions not addressed in the original publication of Lu et al. (2005). Firstly, we performed multiple testing to identify pairs of miRNAs with significantly correlated expression profiles. The fact that a majority of pairwise correlations are significantly different from zero, even after adjusting for multiple tests (nominal FWER level 0.05), suggests a great deal of structure in miRNA expression. This prompted us to perform hierarchical clustering of the miRNA profiles. The HOPACH algorithm yielded sensible ordering of the miRNAs, with groups of similarly expressed miRNAs visualized as blocks in the pseudo-color image of the $J \times J$ ($J = 155$) correlation matrix (Figure 7). In order to focus on the most correlated pairs of miRNAs, the correlation tests could be repeated with a null value larger than zero, e.g., $H_0(j, k) = I(\rho(j, k) \leq 0.5)$, $j = 1, \dots, J$, $k = j + 1, \dots, J$. Further investigation of clusters of co-expressed miRNAs could reveal biologically interesting connections between miRNAs.

We note that the large number of significant findings in both miRNA testing problems is likely due to a reasonably large sample size ($n = 186$) relative to the number of tests ($M = 155$ regression coefficients and $M = 11,935$ correlation coefficients), as compared to similar studies of mRNA expression. Nonetheless, the analysis of a rich dataset using novel and rigorous statistical methods highlights the possibility for meaningful biological and medical discovery from genomic studies.

Software

Our proposed resampling-based multiple testing procedures are implemented in the R package `multtest`, released as part of the Bioconductor Project, an open-source software project for the analysis of biomedical and genomic data (Pollard et al. (2005); `multtest` package, Version 1.6.0, Bioconductor Release 1.6, www.bioconductor.org). Birkner et al. (2005b) illustrate the implementation in SAS (Version 9) of the bootstrap-based single-step maxT procedure and augmentation procedures for controlling the generalized family-wise error rate (gFWER) and tail probabilities for the proportion of false positives (TPFP) among the rejected hypotheses.

The hierarchical ordered partitioning and collapsing hybrid (HOPACH) algorithm is implemented in the Bioconductor R package **hopach** (Pollard and van der Laan (2005); **hopach** package, Version 1.1.1, Bioconductor Release 1.6).

References

- M. D. Birkner, M. Courtine, S. Dudoit, M. J. van der Laan, K. Clément, and J-D. Zucker. Statistical methods for detecting genotype/phenotype associations in the ObeLinks Project. Technical report, Division of Biostatistics, University of California, Berkeley, 2005a. (In preparation).
- M. D. Birkner, K. S. Pollard, M. J. van der Laan, and S. Dudoit. Multiple testing procedures and applications to genomics. Technical Report 168, Division of Biostatistics, University of California, Berkeley, 2005b. URL www.bepress.com/ucbbiostat/paper168.
- M. D. Birkner, S. E. Sinisi, and M. J. van der Laan. Multiple testing and data adaptive regression: An application to hiv-1 sequence data. *Statistical Applications in Genetics and Molecular Biology*, 4(1):Article 8, 2005c. URL www.bepress.com/sagmb/vol4/iss1/art8.
- G. A. Calin, C. Sevignani, C. D. Dumitru, T. Hyslop, E. Noch, S. Yendamuri, M. Shimizu, S. Rattan, F. Bullrich, M. Negrini, and C. M. Croce. Human microRNA genes are frequently located at fragile sites and genomic regions involved in cancers. *Proc. Natl. Acad. Sci.*, 101(9):2999–3004, 2004.
- A. M. Cheng, M. W. Byrom, J. Shelton, and L. P. Ford. Antisense inhibition of human miRNAs and indications for an involvement of miRNA in cell growth and apoptosis. *Nucleic Acids Research*, 33(4):1290–1297, 2005.
- S. Dudoit, J. P. Shaffer, and J. C. Boldrick. Multiple hypothesis testing in microarray experiments. *Statistical Science*, 18(1):71–103, 2003.
- S. Dudoit and M. J. van der Laan. *Multiple Testing Procedures and Applications to Genomics*. Springer, 2005. (In preparation).
- S. Dudoit, M. J. van der Laan, and M. D. Birkner. Multiple testing procedures for controlling tail probability error rates. Technical Report 166,

- Division of Biostatistics, University of California, Berkeley, 2004a. URL www.bepress.com/ucbbiostat/paper166.
- S. Dudoit, M. J. van der Laan, and K. S. Pollard. Multiple testing. Part I. Single-step procedures for control of general Type I error rates. *Statistical Applications in Genetics and Molecular Biology*, 3(1):Article 13, 2004b. URL www.bepress.com/sagmb/vol3/iss1/art13.
- Y. Hochberg and A. C. Tamhane. *Multiple Comparison Procedures*. Probability and Mathematical Statistics. Wiley–Interscience, 1987.
- S. Keleş, M. J. van der Laan, S. Dudoit, and S. E. Cawley. Multiple testing methods for ChIP-Chip high density oligonucleotide array data. Technical Report 147, Division of Biostatistics, University of California, Berkeley, 2004. URL www.bepress.com/ucbbiostat/paper147.
- E. L. Korn, J. F. Troendle, L. M. McShane, and R. Simon. Controlling the number of false discoveries: application to high-dimensional genomic data. *Journal of Statistical Planning and Inference*, 2004. URL linus.nci.nih.gov/~brb/TechReport.htm. (To appear).
- J. Lu, G. Getz, E. A. Miska, E. Alvarez-Saavedra, J. Lamb, D. Peck, A. Sweet-Cordero, B. L. Ebert, R. H. Mak, A. A. Ferrando, J. R. Downing, T. Jacks, H. R. Horvitz, and T. R. Golub. MicroRNA expression profiles classify human cancers. *Nature*, 435(9):834–838, 2005. URL www.broad.mit.edu/cancer/pub/miGCM.
- M. Z. Michael, S. M. O’Connor, N. G. van Holst Pellekaan, G. P. Young, and R. J. James. Reduced accumulation of specific microRNAs in colorectal neoplasia. *Molecular Cancer Research*, 1(12):882–891, 2003.
- K. A. O’Donnell, E. A. Wentzel, K. I. Zeller, C. V. Dang, and J. T. Mendell. c-Myc-regulated microRNAs modulate E2F1 expression. *Nature*, 435(7043):839–843, 2005.
- K. S. Pollard, S. Dudoit, and M. J. van der Laan. *Bioinformatics and Computational Biology Solutions Using R and Bioconductor*, chapter Multiple Testing Procedures: R `multtest` Package and Applications to Genomics. Springer-Verlag, New York, 2005. URL www.bepress.com/ucbbiostat/paper164. (To appear).

- K. S. Pollard and M. J. van der Laan. Choice of a null distribution in resampling-based multiple testing. *Journal of Statistical Planning and Inference*, 125(1–2):85–100, 2004.
- K. S. Pollard and M. J. van der Laan. *Bioinformatics and Computational Biology Solutions Using R and Bioconductor*, chapter Cluster analysis of genomic data with applications in R. Springer-Verlag, New York, 2005. URL www.bepress.com/ucbbiostat/paper167. (To appear).
- J. F. Troendle. A stepwise resampling method of multiple hypothesis testing. *Journal of the American Statistical Association*, 90:370–378, 1995.
- J. F. Troendle. A permutational step-up method of testing multiple outcomes. *Biometrics*, 52:846–859, 1996.
- M. J. van der Laan, M. D. Birkner, and A. E. Hubbard. Resampling based multiple testing procedure controlling tail probability of the proportion of false positives. Technical Report 172, Division of Biostatistics, University of California, Berkeley, 2005. URL www.bepress.com/ucbbiostat/paper172.
- M. J. van der Laan, S. Dudoit, and K. S. Pollard. Augmentation procedures for control of the generalized family-wise error rate and tail probabilities for the proportion of false positives. *Statistical Applications in Genetics and Molecular Biology*, 3(1):Article 15, 2004a. URL www.bepress.com/sagmb/vol3/iss1/art15.
- M. J. van der Laan, S. Dudoit, and K. S. Pollard. Multiple testing. Part II. Step-down procedures for control of the family-wise error rate. *Statistical Applications in Genetics and Molecular Biology*, 3(1):Article 14, 2004b. URL www.bepress.com/sagmb/vol3/iss1/art14.
- M. J. van der Laan and K. S. Pollard. Hybrid clustering of gene expression data with visualization. *Journal of Statistical Planning and Inference*, 117: 275–303, 2003.
- P. H. Westfall and S. S. Young. *Resampling-based multiple testing: Examples and methods for p-value adjustment*. John Wiley & Sons, 1993.

Table 1: *Type I and Type II errors in multiple hypothesis testing.* This table summarizes the different types of decisions and errors in multiple hypothesis testing. The number of rejected hypotheses is $R_n \equiv |\mathcal{R}_n| = \sum_{m=1}^M \mathbf{I}(T_n(m) \in \mathcal{C}_n(m))$, the number of Type I errors is $V_n \equiv |\mathcal{R}_n \cap \mathcal{H}_0| = \sum_{m \in \mathcal{H}_0} \mathbf{I}(T_n(m) \in \mathcal{C}_n(m))$, and the number of Type II errors is $U_n \equiv |\mathcal{R}_n^c \cap \mathcal{H}_1| = \sum_{m \in \mathcal{H}_1} \mathbf{I}(T_n(m) \notin \mathcal{C}_n(m))$.

		Null hypotheses		
		not rejected	rejected	
Null hypotheses	true	$ \mathcal{R}_n^c \cap \mathcal{H}_0 $	$V_n = \mathcal{R}_n \cap \mathcal{H}_0 $ (Type I)	$h_0 = \mathcal{H}_0 $
	false	$U_n = \mathcal{R}_n^c \cap \mathcal{H}_1 $ (Type II)	$ \mathcal{R}_n \cap \mathcal{H}_1 $	$h_1 = \mathcal{H}_1 $
		$M - R_n$	$R_n = \mathcal{R}_n $	M

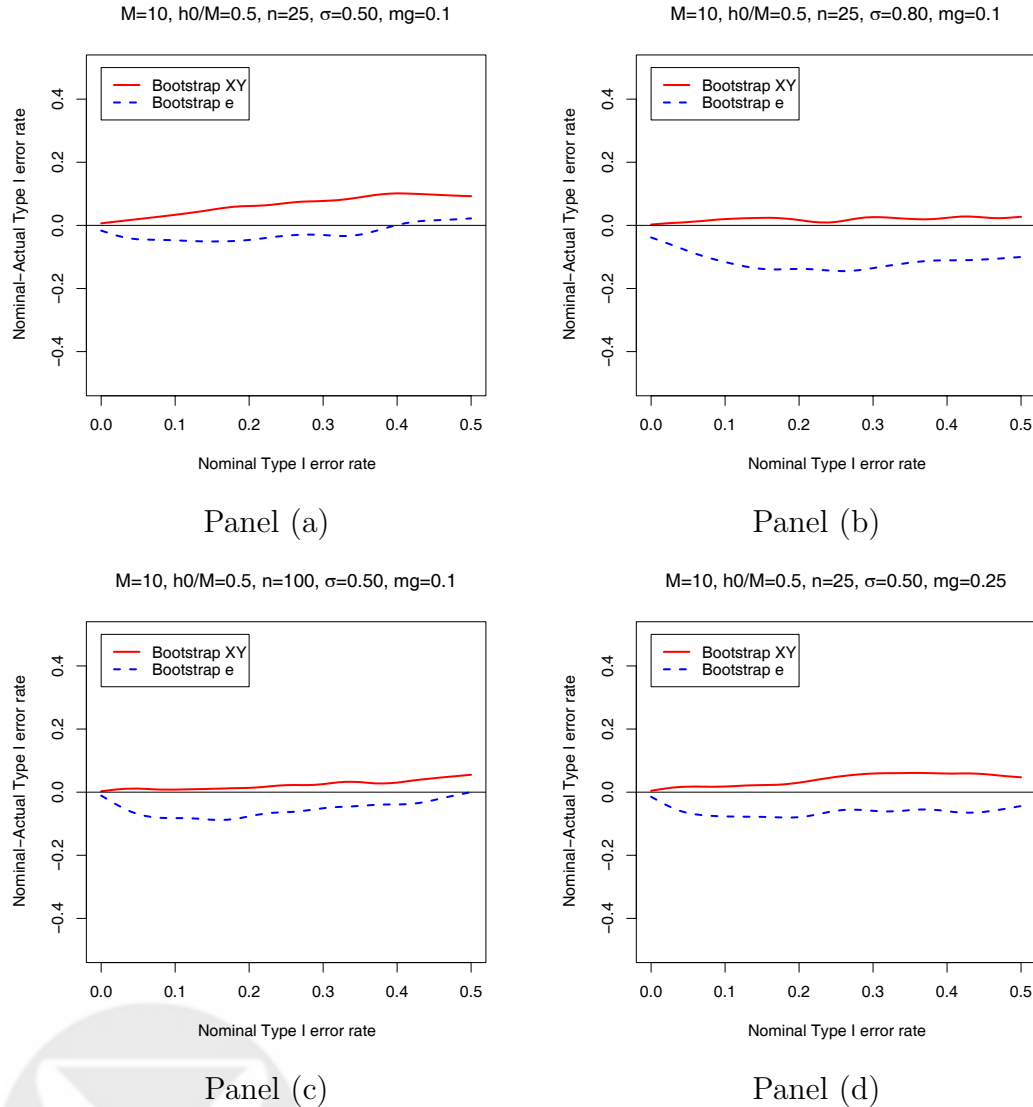


Figure 1: *Simulation Study 1: Tests of linear regression coefficients, Type I error control comparison.* Plots of differences between nominal and actual Type I error rates vs. nominal Type I error rate, for single-step maxT procedure based on general non-parametric bootstrap null distribution **Bootstrap XY** and residual bootstrap null distribution **Bootstrap e**. The null hypotheses are tested using the t -statistics of Equation (18). Panel (a): Model with sample size $n = 25$; $M = 10$ null hypotheses; common covariance for the covariates $\varsigma = 0.50$; proportion $h_0/M = 0.50$ of true null hypotheses; magnitude parameter for alternative regression coefficients $\mu = 0.10$. Panel (b): $n = 25$; $M = 10$; $\varsigma = 0.80$; $h_0/M = 0.50$; $\mu = 0.10$. Panel (c): $n = 100$; $M = 10$; $\varsigma = 0.50$; $h_0/M = 0.50$; $\mu = 0.10$. Panel (d): $n = 25$; $M = 10$; $\varsigma = 0.50$; $h_0/M = 0.50$; $\mu = 0.25$.

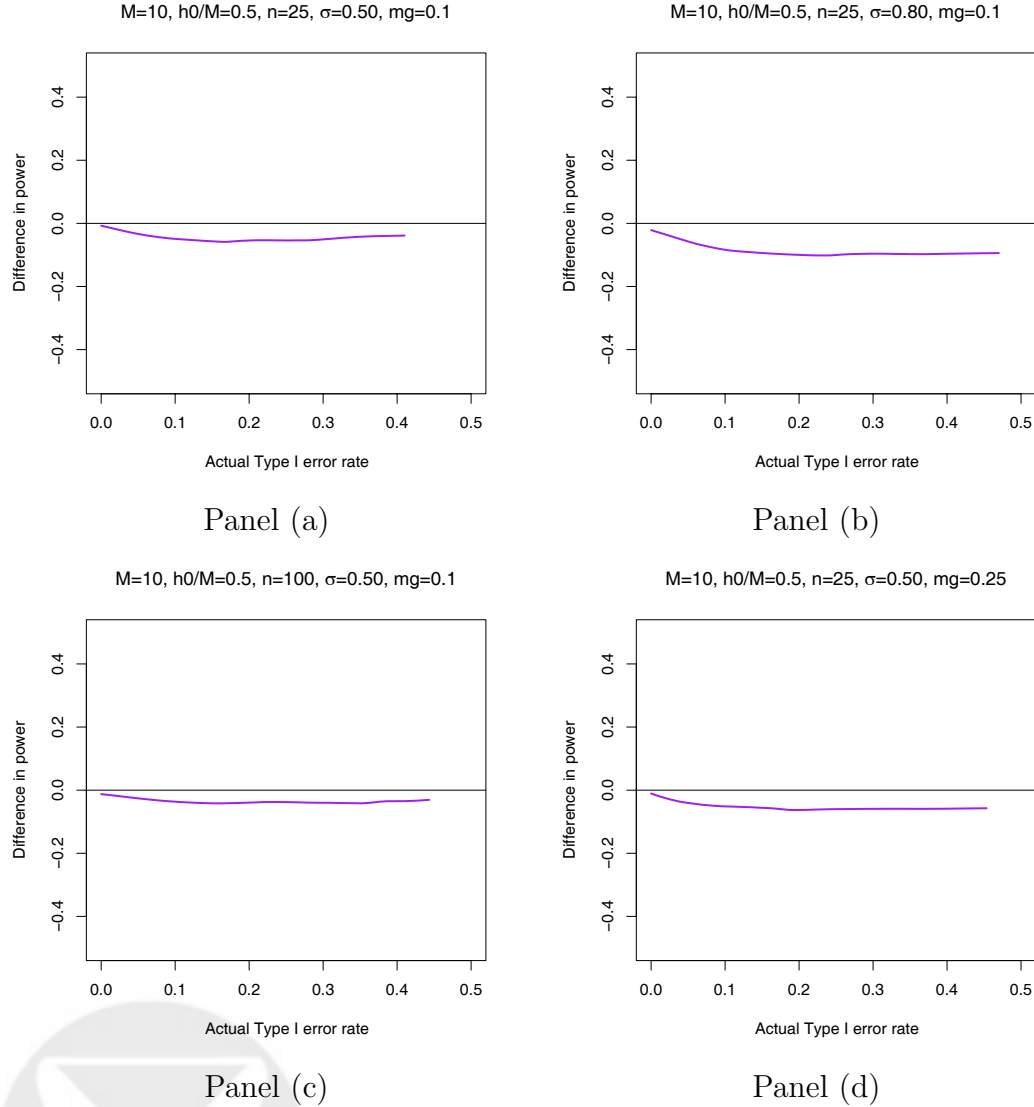


Figure 2: *Simulation Study 1: Tests of linear regression coefficients, power comparison.* Plots of difference in power vs. actual Type I error rate, for single-step maxT procedure based on general non-parametric bootstrap null distribution **Bootstrap XY** and residual bootstrap null distribution **Bootstrap e**. The null hypotheses are tested using the t -statistics of Equation (18). Positive differences indicate greater power for **Bootstrap XY**. Panel (a): Model with sample size $n = 25$; $M = 10$ null hypotheses; common covariance for the covariates $\varsigma = 0.50$; proportion $h_0/M = 0.50$ of true null hypotheses; magnitude parameter for alternative regression coefficients $\mu = 0.10$. Panel (b): $n = 25$; $M = 10$; $\varsigma = 0.80$; $h_0/M = 0.50$; $\mu = 0.10$. Panel (c): $n = 100$; $M = 10$; $\varsigma = 0.50$; $h_0/M = 0.50$; $\mu = 0.10$. Panel (d): $n = 25$; $M = 10$; $\varsigma = 0.50$; $h_0/M = 0.50$; $\mu = 0.25$.

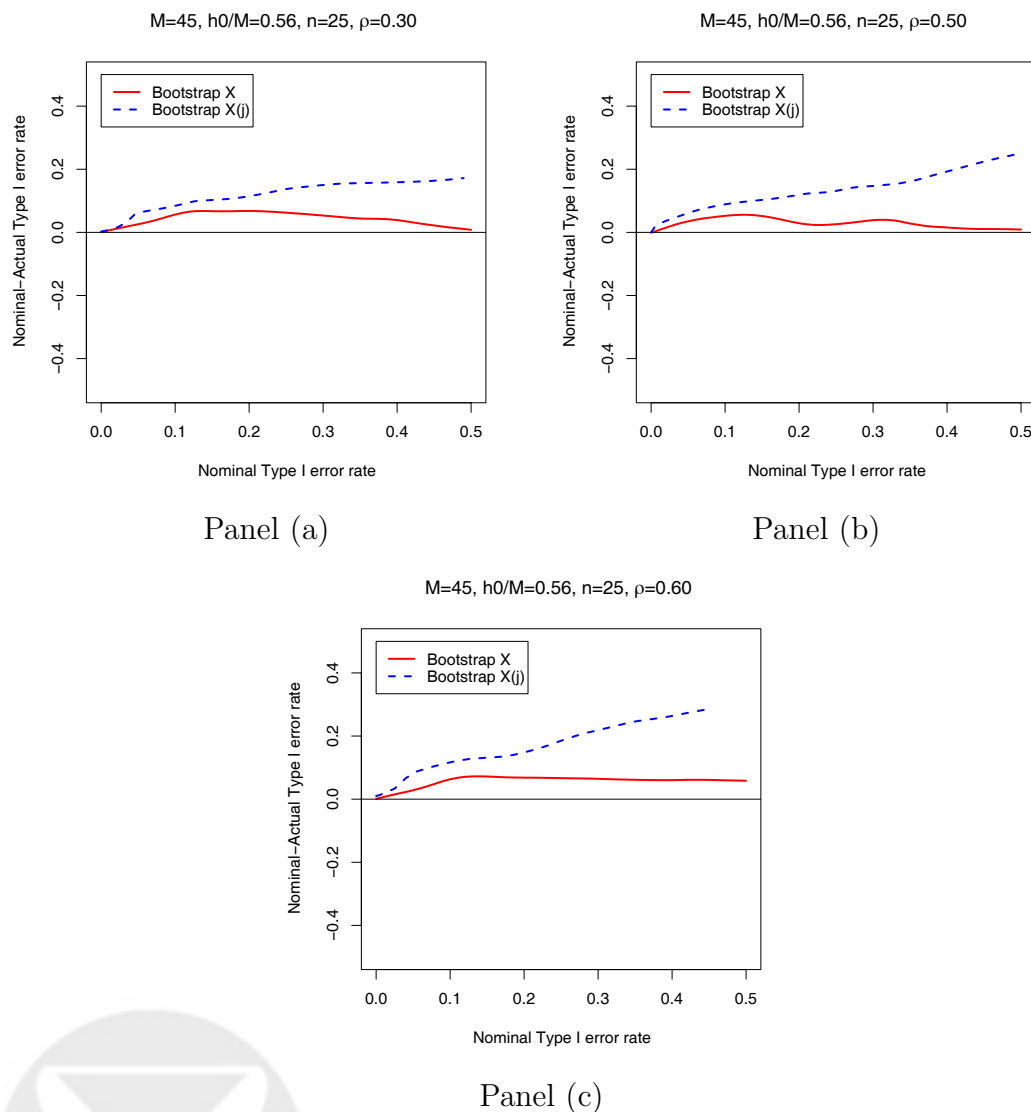


Figure 3: *Simulation Study 2: Tests of correlation coefficients, Type I error control comparison.* Plots of differences between nominal and actual Type I error rates vs. nominal Type I error rate, for single-step maxT procedure based general non-parametric bootstrap null distribution **Bootstrap X** and independent covariates bootstrap null distribution **Bootstrap X(j)**. The null hypotheses are tested using the t -statistics of Equation (24). Model with sample size $n = 25$; $M = 45$ null hypotheses; proportion $h_0/M = 25/45$ of true null hypotheses. Panel (a): common correlation coefficient for the two blocks $\rho = 0.30$. Panel (b): $\rho = 0.50$. Panel (c): $\rho = 0.60$.

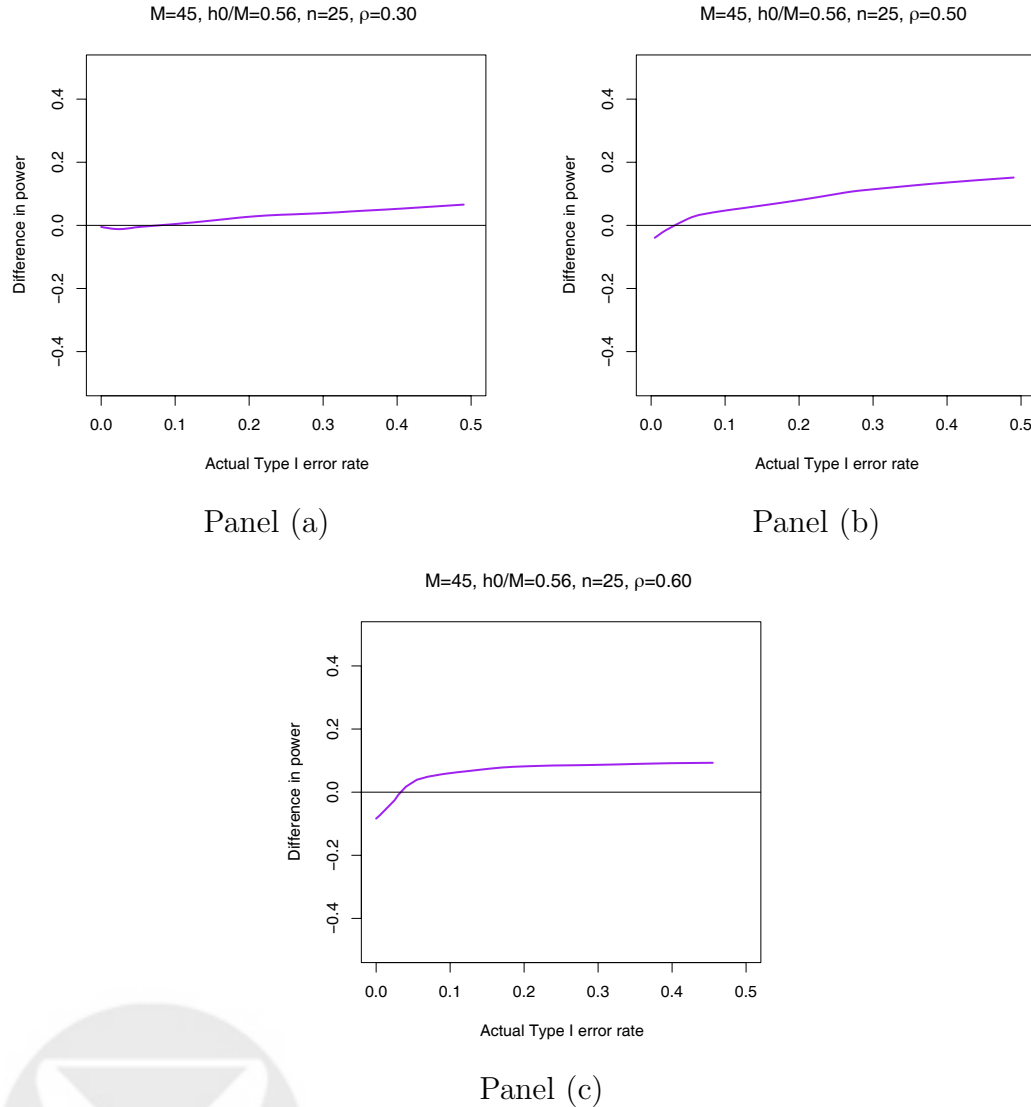


Figure 4: *Simulation Study 2: Tests of correlation coefficients, power comparison.* Plots of difference in power vs. actual Type I error rate, for single-step maxT procedure based on general non-parametric bootstrap null distribution **Bootstrap X** and independent covariates bootstrap null distribution **Bootstrap X(j)**. The null hypotheses are tested using the t -statistics of Equation (24). Positive differences indicate greater power for **Bootstrap X**. Model with sample size $n = 25$; $M = 45$ null hypotheses; proportion $h_0/M = 25/45$ of true null hypotheses. Panel (a): common correlation coefficient for the two blocks $\rho = 0.30$. Panel (b): $\rho = 0.50$. Panel (c): $\rho = 0.60$.

Table 2: *miRNA data analysis: Tests of logistic regression coefficients.* The table reports the names, target sequences, and test statistics, for the 36 miRNAs most significantly differentially expressed between cancerous and healthy tissues, according to bootstrap-based single-step maxT Procedure 3. All 36 miRNAs have adjusted p -values approximately equal to zero. miRNAs are ordered by test statistic $T_n(j)$, with most significant (largest in absolute value) first. Negative test statistics indicate under-expression in cancerous compared to healthy tissues. The target sequence is the reverse complement of the miRNA sequence, which identifies potential binding sites for the miRNA.

Located in minimal deleted regions, minimal amplified regions, and break-point regions involved in human cancers (Calin et al., 2004) .

Name	miRNA target sequence	Test statistic
hsa-miR-98	UGAGGUAGUAAGUUGUAUUGUU	-4.88
hsa-miR-28	AAGGAGCUCACAGUCUAUUGAG	-4.79
hsa-miR-196	UAGGUAGUUUCAUGUUGUUGG	-4.79
hsa-miR-30a	CUUUCAGUCGGAUGUUUGCAGC	-4.78
hsa-miR-30e	UGUAAACAUCUUGACUGGA	-4.78
hsa-miR-99a#	AACCCGUAGAUCGUAUCUUGUG	-4.77
rno-miR-335	UCAAGAGCAAUAAACGAAAAAUGU	-4.72
hsa-let-7e	UGAGGUAGGAGGUUGUAUAGU	-4.69
hsa-miR-23b#	AUCACAUUGCCAGGGAUUACCAC	-4.67
hsa-miR-99b	CACCCGUAGAACCGACCUUGCG	-4.67
hsa-miR-214	ACAGCAGGCACAGACAGGCAG	-4.67
hsa-miR-30b	UGUAAACAUCUACACUCAGC	-4.66
hsa-miR-30c	UGUAAACAUCUACACUCUCAGC	-4.66
mmu-miR-338	UCCAGCAUCAGUGAUUUUGUUGA	-4.65
hsa-miR-103	AGCAGCAUUGUACAGGGCUAUGA	-4.64
hsa-miR-185	UGGAGAGAAAGGCAGUUC	-4.63
rno-miR-151*	UCGAGGAGCUCACAGUUCUAGUA	-4.62
hsa-miR-20_(sub_1)	UAAAGUGCUUAUAGUGCAGGUAG	-4.61
hsa-miR-100#	AACCCGUAGAUCGUAUCUUGUG	-4.61
hsa-miR-22#	AAGCUGCCAGUUGAAGAACUGU	-4.60
rno-miR-129*	AAGCCCUUACCCCAAAAAGCAU	-4.60
hsa-let-7d#	AGAGGUAGUAGGUUGCAUAGU	-4.58
hsa-miR-107	AGCAGCAUUGUACAGGGCUAUCA	-4.58
rno-miR-352	AGAGUAGUAGGUUGCAUAGUA	-4.58
hsa-miR-32	UAUUGCACAUAUACUAAGUUGC	-4.57
hsa-miR-197	UUCACCAACCUUCUCCACCCAGC	-4.57
mmu-miR-342	UCUCACACAGAAUUCGACCCCGUC	-4.56
hsa-miR-128b	UCACAGUGAACCGGUCUCUUUC	-4.51
mmu-miR-324-5p	CGCAUCCCCUAGGGCAUUGGUGU	-4.51
hsa-miR-126*	CAUUAUUACUUUUGGUACGCG	-4.50
hsa-miR-19b	UGUGCAAUCCAUUGCAAAACUGA	-4.49
mmu-miR-151_(sub_1)	ACUAGACUGAGGCUCUUGAGG	-4.49
hsa-let-7i	UGAGGUAGUAGUUUGUGCU	-4.48
hsa-miR-199a*	UACAGUAGUCUGCACAUUGGUU	-4.48
hsa-miR-10b	UACCCUGUAGAACCGAAUUUGU	-4.47
mmu-miR-292-3p	AAGUGCCGCCAGGUUUUGAGUGU	-4.46

Table 3: *miRNA data analysis: Tests of correlation coefficients.* The table reports the names and correlation coefficients for the twenty pairs of miRNAs with the most significantly correlated expression profiles, according to bootstrap-based single-step maxT Procedure 3. miRNAs are ordered by test statistic $T_n(m)$, with most significant (largest in absolute value) first. Several pairs are composed of miRNAs in the same family, e.g., hsa-miR-10a and hsa-miR-10b.

Up-regulated by the proto-oncogene c-MYC (O'Donnell et al., 2005). † Increases cell growth in lung carcinomas (Cheng et al., 2005). ‡ Expressed at lower levels in cancerous and pre-cancerous compared to normal colon tissue (Michael et al., 2003).

Names		Correlation coefficient
hsa-miR-106a#	hsa-miR-17-5p#	0.99
mmu-miR-200b	hsa-miR-200b	0.99
mmu-miR-200b	hsa-miR-200c	0.99
hsa-miR-107†	hsa-miR-103	0.99
hsa-miR-200b	hsa-miR-200c	0.99
hsa-miR-145‡	hsa-miR-143‡	0.98
hsa-miR-199a_(sub_1)	mmu-miR-199b	0.98
hsa-miR-17-5p	hsa-miR-20_(sub_1)	0.97
hsa-miR-19a#	hsa-miR-19b#	0.97
hsa-miR-29a	hsa-miR-30a*	0.97
hsa-miR-181a	hsa-miR-181c	0.97
hsa-miR-199a_(sub_1)	hsa-miR-199a*	0.97
hsa-miR-29b_(sub_2)	hsa-miR-29c	0.97
hsa-miR-199a*	mmu-miR-199b	0.96
hsa-miR-200a	hsa-miR-141	0.96
hsa-miR-20_(sub_1)#	mmu-miR-106a	0.96
hsa-miR-106a	hsa-miR-20_(sub_1)#	0.96
hsa-miR-200a	hsa-miR-200a	0.96
hsa-miR-23b	hsa-miR-23a	0.96
hsa-miR-10a	hsa-miR-10b	0.96

Logistic Regression: Cancer~miRNA+tissue

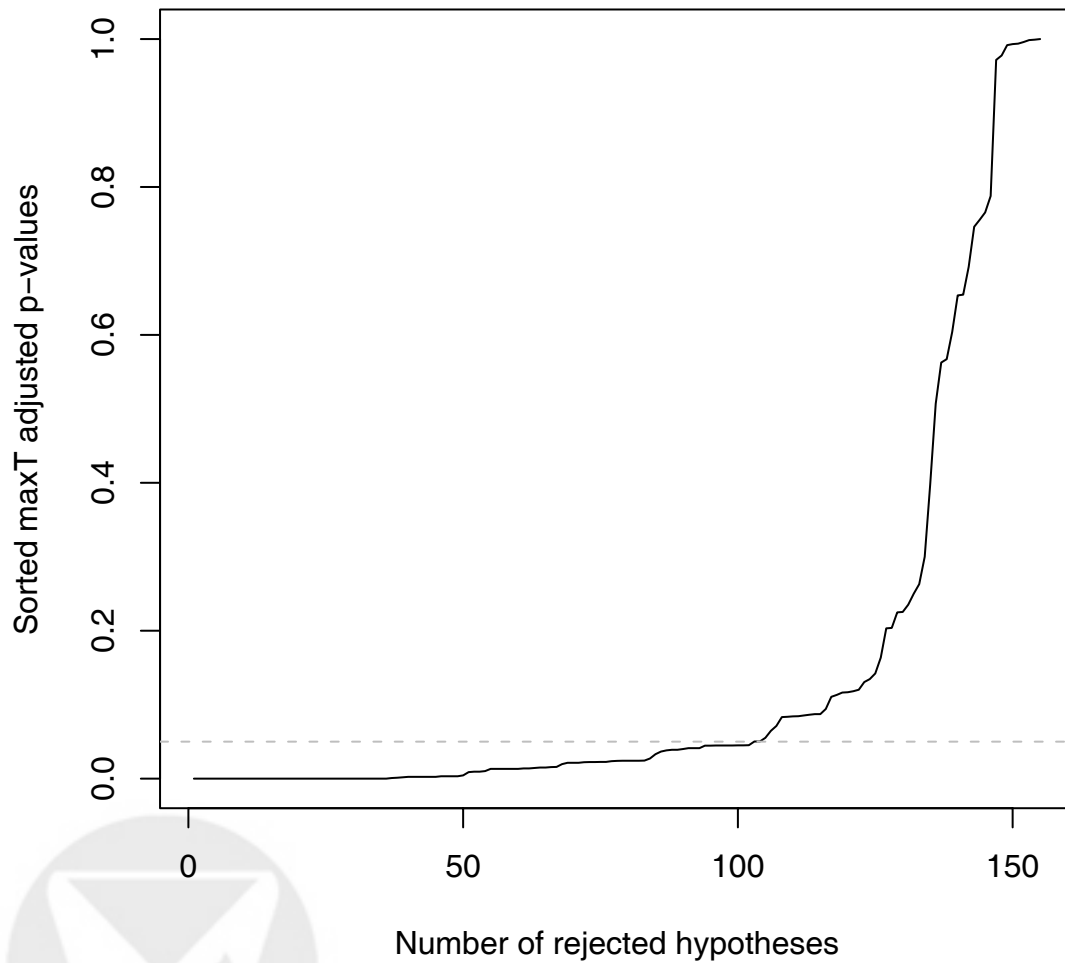


Figure 5: *miRNA data analysis: Tests of logistic regression coefficients.* Plot of sorted adjusted p -values for bootstrap-based single-step maxT Procedure 3.

Pairwise Correlations between miRNAs

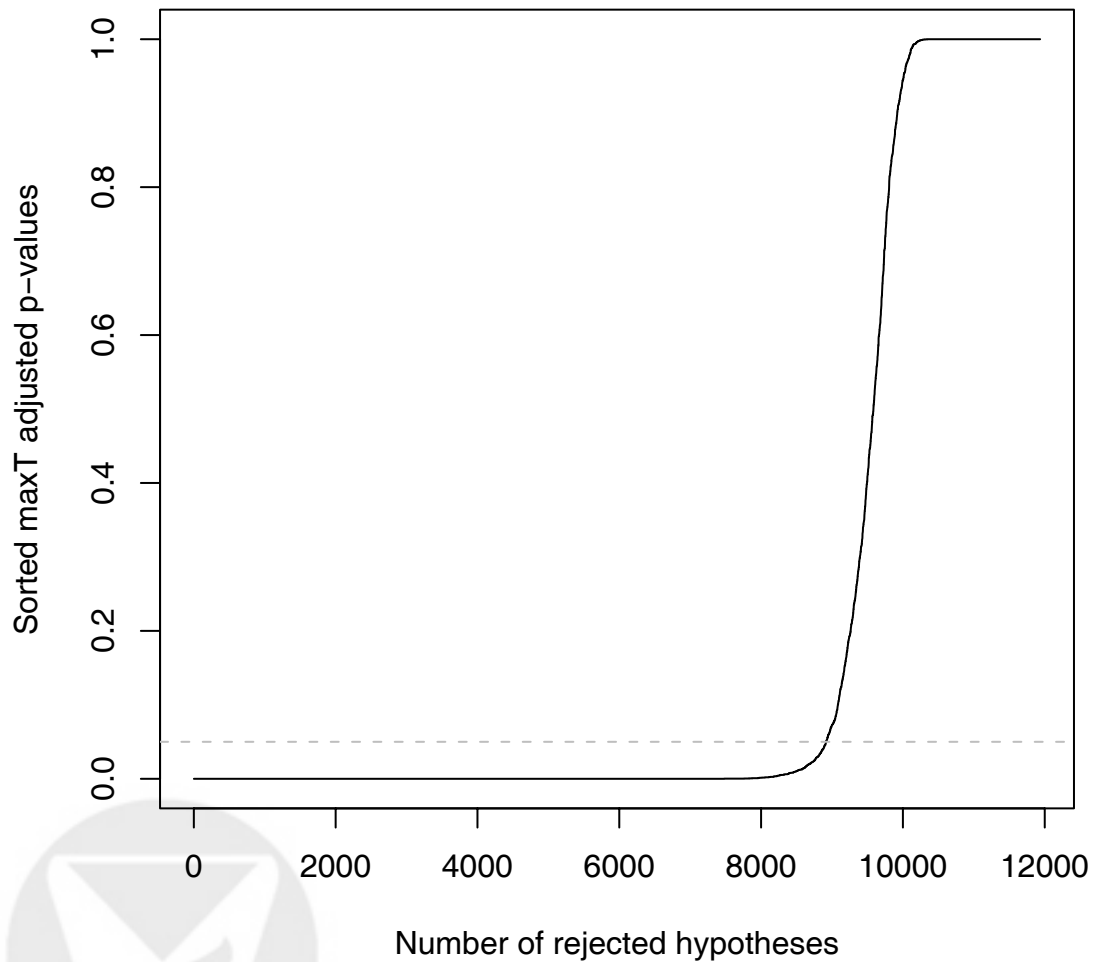


Figure 6: *miRNA data analysis: Tests of correlation coefficients.* Plot of sorted adjusted p -values for bootstrap-based single-step maxT Procedure 3.

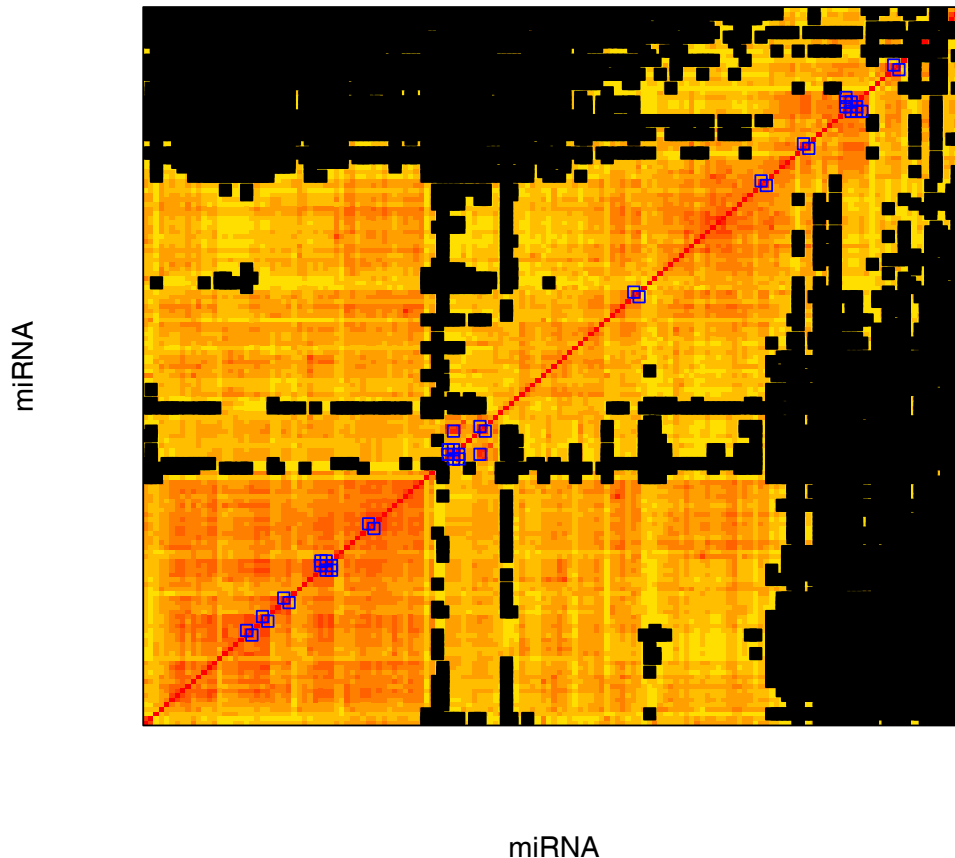


Figure 7: *miRNA data analysis: HOPACH clustering of miRNA expression profiles using correlation distance.* The figure provides a pseudo-color image of the 155×155 correlation matrix for the expression profiles of the $J = 155$ miRNAs, with rows and columns ordered according to the final level of the hierarchical tree of miRNA clusters. Pairwise correlations not significantly different from zero are displayed in black. Remaining correlations are represented using a white (anti-correlated) to red (positively-correlated) color palette. Groups of co-expressed miRNAs appear as red blocks on the diagonal of the correlation matrix. The twenty most significantly correlated pairs of miRNAs from Table 3 are marked in blue.